

Alípio Jorge Luís Torgo
Pavel Brazdil Rui Camacho
João Gama (Eds.)

LNAI 3721

Knowledge Discovery in Databases: PKDD 2005

9th European Conference on Principles and Practice
of Knowledge Discovery in Databases
Porto, Portugal, October 2005, Proceedings



 Springer

Lecture Notes in Artificial Intelligence 3721

Edited by J. G. Carbonell and J. Siekmann

Subseries of Lecture Notes in Computer Science

Alípio Jorge Luís Torgo
Pavel Brazdil Rui Camacho
João Gama (Eds.)

Knowledge Discovery in Databases: PKDD 2005

9th European Conference on Principles and Practice
of Knowledge Discovery in Databases
Porto, Portugal, October 3-7, 2005
Proceedings



Springer

Series Editors

Jaime G. Carbonell, Carnegie Mellon University, Pittsburgh, PA, USA
Jörg Siekmann, University of Saarland, Saarbrücken, Germany

Volume Editors

Alípio Jorge
Luís Torgo
Pavel Brazdil
João Gama
LIACC/FEP, University of Porto
Rua de Ceuta, 118, 6°, 4050-190 Porto, Portugal
E-mail: {amjorge,ltorgo,pbrazdil,jgama}@liacc.up.pt

Rui Camacho
LIACC/FEUP, University of Porto
Rua de Ceuta, 118, 6°, 4050-190 Porto, Portugal
E-mail: rcamacho@fe.up.pt

Library of Congress Control Number: 2005933047

CR Subject Classification (1998): I.2, H.2, J.1, H.3, G.3, I.7, F.4.1

ISSN 0302-9743
ISBN-10 3-540-29244-6 Springer Berlin Heidelberg New York
ISBN-13 978-3-540-29244-9 Springer Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable to prosecution under the German Copyright Law.

Springer is a part of Springer Science+Business Media

springeronline.com

© Springer-Verlag Berlin Heidelberg 2005
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India
Printed on acid-free paper SPIN: 11564126 06/3142 5 4 3 2 1 0

Preface

The European Conference on Machine Learning (ECML) and the European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD) were jointly organized this year for the fifth time in a row, after some years of mutual independence before. After Freiburg (2001), Helsinki (2002), Cavtat (2003) and Pisa (2004), Porto received the 16th edition of ECML and the 9th PKDD in October 3–7.

Having the two conferences together seems to be working well: 585 different paper submissions were received for both events, which maintains the high submission standard of last year. Of these, 335 were submitted to ECML only, 220 to PKDD only and 30 to both. Such a high volume of scientific work required a tremendous effort from Area Chairs, Program Committee members and some additional reviewers. On average, PC members had 10 papers to evaluate, and Area Chairs had 25 papers to decide upon. We managed to have 3 highly qualified independent reviews per paper (with very few exceptions) and one additional overall input from one of the Area Chairs. After the authors' responses and the online discussions for many of the papers, we arrived at the final selection of 40 regular papers for ECML and 35 for PKDD. Besides these, 32 others were accepted as short papers for ECML and 35 for PKDD. This represents a joint acceptance rate of around 13% for regular papers and 25% overall. We thank all involved for all the effort with reviewing and selection of papers.

Besides the core technical program, ECML and PKDD had 6 invited speakers, 10 workshops, 8 tutorials and a Knowledge Discovery Challenge. Our special thanks to the organizers of the individual workshops and tutorials and to the workshop and tutorial chairs Floriana Esposito and Dunja Mladenić and to the challenge organizer Petr Berka. A very special word to Richard van de Stadt for all his competence and professionalism in the management of CyberChairPRO. Our thanks also to everyone from the Organization Committee mentioned further on who helped us with the organization. Our acknowledgement also to Rodolfo Matos and Assunção Costa Lima for providing logistic support.

Our acknowledgements to all the sponsors, Fundação para a Ciência e Tecnologia (FCT), LIACC-NIAAD, Faculdade de Engenharia do Porto, Faculdade de Economia do Porto, KDubiq – Knowledge Discovery in Ubiquitous Environments—Coordinated Action of FP6, Salford Systems, Pascal Network of Excellence, PSE/SPSS, ECCAI and Comissão de Viticultura da Região dos Vinhos Verdes. We also wish to express our gratitude to all other individuals and institutions not explicitly mentioned in this text who somehow contributed to the success of these events.

Finally, our word of appreciation to all the authors who submitted papers to the main conferences and their workshops, without whom none of this would have been possible.

July 2005

Alípio Jorge, Luís Torgo, Pavel Brazdil,
Rui Camacho and João Gama

Organization

ECML/PKDD 2005 Organization

Executive Committee

General Chair

Pavel Brazdil (LIACC/FEP, Portugal)

Program Chairs

Rui Camacho (LIACC/FEUP, Portugal)

João Gama (LIACC/FEP, Portugal)

Alípio Jorge (LIACC/FEP, Portugal)

Luís Torgo (LIACC/FEP, Portugal)

Workshop Chair

Floriana Esposito (University of Bari, Italy)

Tutorial Chair

Dunja Mladenić (Jozef Stefan Institute, Slovenia)

Challenge Chairs

Petr Berka (University of Economics, Czech Republic)

Bruno Crémilleux (Université de Caen, France)

Local Organization Committee

Pavel Brazdil, Alípio Jorge, Rui Camacho, Luís Torgo and João Gama, with the help of people from LIACC-NIAAD, University of Porto, Portugal, Rodolfo Matos, Pedro Quelhas Brito, Fabrice Colas, Carlos Soares, Pedro Campos, Rui Leite, Mário Amado Alves, Pedro Rodrigues; and from IST, Lisbon, Portugal, Cláudia Antunes.

Steering Committee

Nada Lavrač, Jozef Stefan Institute, Slovenia

Dragan Gamberger, Rudjer Boskovic Institute, Croatia

Ljupčo Todorovski, Jozef Stefan Institute, Slovenia

Hendrik Blockeel, Katholieke Universiteit Leuven, Belgium

Tapio Elomaa, Tampere University of Technology, Finland

Heikki Mannila, Helsinki Institute for Information Technology, Finland

Hannu T.T. Toivonen, University of Helsinki, Finland
Jean-François Boulicaut, INSA-Lyon, France
Floriana Esposito, University of Bari, Italy
Fosca Giannotti, ISTI-CNR, Pisa, Italy
Dino Pedreschi, University of Pisa, Italy

Area Chairs

Michael R. Berthold, Germany	Alípio Jorge, Portugal
Elisa Bertino, Italy	Hillol Kargupta, USA
Ivan Bratko, Slovenia	Pedro Larranaga, Spain
Pavel Brazdil, Portugal	Ramon López de Mántaras, Spain
Carla E. Brodley, USA	Dunja Mladenić, Slovenia
Rui Camacho, Portugal	Hiroshi Motoda, Japan
Luc Dehaspe, Belgium	José Carlos Príncipe, USA
Peter Flach, UK	Tobias Scheffer, Germany
Johannes Fürnkranz, Germany	Michele Sebag, France
João Gama, Portugal	Peter Stone, USA
Howard J. Hamilton, Canada	Luís Torgo, Portugal
Thorsten Joachims, USA	Gerhard Widmer, Austria

Program Committee

Jesus Aguilar, Spain	Susan Craw, Scotland
Paulo Azevedo, Portugal	Bruno Cremilleux, France
Michael Bain, Australia	James Cussens, UK
José Luís Balcázar, Spain	Luc Dehaspe, Belgium
Elena Baralis, Italy	Vasant Dhar, USA
Bettina Berendt, Germany	Sašo Džeroski, Slovenia
Petr Berka, Czech Republic	Tapio Elomaa, Finland
Michael R. Berthold, Germany	Floriana Esposito, Italy
Elisa Bertino, Italy	Jianping Fan, USA
Hendrik Blockeel, Belgium	Ronen Feldman, Israel
José Luís Borges, Portugal	Cèsar Ferri, Spain
Francesco Bonchi, Italy	Peter Flach, UK
Henrik Boström, Sweden	Eibe Frank, New Zealand
Marco Botta, Italy	Alex Freitas, UK
Jean-François Boulicaut, France	Johannes Fürnkranz, Germany
Pavel Brazdil, Portugal	Thomas Gabriel, Germany
Carla E. Brodley, USA	João Gama, Portugal
Wray Buntine, Finland	Dragan Gamberger, Croatia
Rui Camacho, Portugal	Minos N. Garofalakis, USA
Amilcar Cardoso, Portugal	Bart Goethals, Belgium
Barbara Catania, Italy	Gunter Grieser, Germany
Jesús Cerquides, Spain	Marko Grobelnik, Slovenia

- Howard J. Hamilton, Canada
 Colin de la Higuera, France
 Melanie Hilario, Switzerland
 Haym Hirsh, USA
 Frank Höppner, Germany
 Andreas Hotho, Germany
 Eyke Hullermeier, German
 Nitin Indurkha, Australia
 Thorsten Joachims, USA
 Aleks Jakulin, Slovenia
 Alípio Jorge, Portugal
 Murat Kantarcioglu, USA
 Hillol Kargupta, USA
 Samuel Kaski, Finland
 Mehmet Kaya, Turkey
 Eamonn Keogh, USA
 Roni Khardon, USA
 Jorg-Uwe Kietz, Switzerland
 Joerg Kindermann, Germany
 Ross D. King, UK
 Joost N. Kok, The Netherlands
 Igor Kononenko, Slovenia
 Stefan Kramer, Germany
 Marzena Kryszkiewicz, Poland
 Miroslav Kubat, USA
 Nicholas Kushmerick, Ireland
 Pedro Larranaga, Spain
 Nada Lavrač, Slovenia
 Jure Leskovec, USA
 Jinyan Li, Singapore
 Charles Ling, Canada
 Donato Malerba, Italy
 Giuseppe Manco, Italy
 Ramon López de Mántaras, Spain
 Stan Matwin, Canada
 Michael May, Germany
 Rosa Meo, Italy
 Dunja Mladenić, Slovenia
 Eduardo Morales, Mexico
 Katharina Morik, Germany
 Shinichi Morishita, Japan
 Hiroshi Motoda, Japan
 Steve Moyle, England
 Gholamreza Nakhaeizadeh, Germany
 Claire Nédellec, France
 Richard Nock, France
 Andreas Nürnberger, Germany
 Masayuki Numao, Japan
 Arlindo Oliveira, Portugal
 Georgios Paliouras, Greece
 Dino Pedreschi, Italy
 Johann Petrak, Austria
 Bernhard Pfahringer, New Zealand
 Lubomir Popelinsky, Czech
 José Carlos Príncipe, USA
 Jan Ramon, Belgium
 Zbigniew W. Ras, USA
 Jan Rauch, Czech Republic
 Christophe Rigotti, France
 Gilbert Ritschard, Switzerland
 John Roddick, Australia
 Marques de Sá, Portugal
 Lorenza Saitta, Italy
 Daniel Sanchez, Spain
 Yücel Saygin, Turkey
 Tobias Scheffer, Germany
 Michele Sebag, France
 Marc Sebban, France
 Giovanni Semeraro, Italy
 Arno Siebes, The Netherlands
 Andrzej Skowron, Poland
 Carlos Soares, Portugal
 Maarten van Someren, The Netherlands
 Myra Spiliopoulou, Germany
 Nicolas Spyrtas, France
 Ashwin Srinivasan, India
 Olga Stepankova, Czech Republic
 Reinhard Stolle, Germany
 Peter Stone, USA
 Gerd Stumme, Germany
 Einoshin Suzuki, Japan
 Domenico Talia, Italy
 Ah-Hwee Tan, Singapore
 Evimaria Terzi, Finland
 Ljupčo Todorovski, Slovenia
 Luís Torgo, Portugal
 Shusaku Tsumoto, Japan
 Peter Turney, Canada
 Jaideep Vaidya, USA
 Maria Amparo Vila, Spain

Dietrich Wettschereck, UK
Gerhard Widmer, Austria
Rüdiger Wirth, Germany
Mohammed J. Zaki, USA

Carlo Zaniolo, USA
Djamel A. Zighed, France

Additional Reviewers

Markus Ackermann
Anastasia Analyti
Fabrizio Angiulli
Orlando Anunciação
Giacometti Arnaud
Stella Asimwe
Maurizio Atzori
Korinna Bade
Miriam Baglioni
Juergen Beringer
Philippe Bessières
Matjaž Bevk
Abhilasha Bhargav
Steffen Bickel
Stefano Bistarelli
Jan Blafák
Axel Blumenstock
Janez Brank
Ulf Brefeld
Klaus Brinker
Michael Brückner
Robert D. Burbidge
Toon Calders
Michelangelo Ceci
Tania Cerquitelli
Eugenio Cesario
Vineet Chaoji
Silvia Chiusano
Fang Chu
Chris Clifton
Luís Coelho
Carmela Comito
Gianni Costa
Juan-Carlos Cubero
Agnieszka Dardzinska
Damjan Demšar
Nele Dexters

Christian Diekmann
Isabel Drost
Nicolas Durand
Mohamed Elfeky
Timm Euler
Nicola Fanizzi
Pedro Gabriel Ferreira
Francisco Ferrer
Daan Fierens
Sergio Flesca
Francesco Folino
Blaž Fortuna
Kenichi Fukui
Feng Gao
Yuli Gao
Paolo Garza
Aristides Gionis
Paulo Gomes
Gianluigi Greco
Fabian Güiza
Amaury Habrard
Hakim Hacid
Mohand-Said Hacid
Niina Haiminen
Mark Hall
Christoph Heinz
José Hernández-Orallo
Jochen Hipp
Juan F. Huete
Ali Inan
Ingo Mierswa
Stephanie Jacquemont
François Jacquenet
Aleks Jakulin
Tao Jiang
Wei Jiang
Xing Jiang

Sachindra Joshi
Pierre-Emmanuel Jouve
Michael Steinbach
George Karypis
Steffen Kempe
Arto Klami
Christian Kolbe
Stasinios Konstantopoulos
Matjaz Kukar
Minseok Kwon
Lotfi Lakhhal
Carsten Lanquillon
Dominique Laurent
Yan-Nei Law
Roberto Legaspi
Aurélien Lemaire
Claire Leschi
Jure Leskovec
Francesca A. Lisi
Antonio Locane
Corrado Loglisci
Claudio Lucchese
Sara C. Madeira
Alain-Pierre Manine
M.J. Martin-Bautista
Stewart Massie
Cyrille Masson
Carlo Mastroianni
Nicolas Meger
Carlo Meghini
Taneli Mielikäinen
Mummoorthy Murugesan
Mirco Nanni
Mehmet Ercan Nergiz
Siegfried Nijssen
Janne Nikkilä
Blaž Novak
Lenka Novakova
Merja Oja
Riccardo Ortale
Ignazio Palmisano
Raffaele Perego
Sergios Petridis
Viet Phan-Luong
Frederic Piat
Dimitrios Pierrakos
Ulrich Rückert
María José Ramírez
Ganesh Ramakrishnan
Domenico Redavid
Chiara Renso
Lothar Richter
François Rioult
Céline Robardet
Pedro Rodrigues
Domingo S. Rodriguez-Baena
Luka Šajn
Saeed Salem
D. Sanchez
Eerika Savia
Karlton Sequeira
J.M. Serrano
Shengli Sheng
Javed Siddique
Georgios Sigletos
Fabrizio Silvestri
Arnaud Soulet
Hendrik Stange
Jan Struyf
Henri-Maxime Suchier
Andrea Tagarelli
Julien Thomas
Panayiotis Tsaparas
Yannis Tzitzikas
Jarkko Venna
Celine Vens
Nirmalie Wiratunga
Ghim-Eng Yap
Justin Zhan
Lizhuang Zhao
Igor Zwir

ECML/PKDD 2005 Tutorials

Ontology Learning from Text

Learning Automata as a Basis for Multi-agent Reinforcement Learning

Web Mining for Web Personalization

A Practical Time-Series Tutorial with MATLAB

Mining the Volatile Web

Spectral Clustering

Bioinspired Machine Learning Techniques

Probabilistic Inductive Logic Programming

ECML/PKDD 2005 Workshops

Sub-symbolic Paradigms for Learning in Structured Domains

European Web Mining Forum 2005 (EWMF 2005)

Knowledge Discovery in Inductive Databases (KDID 2005)

Mining Spatio-temporal Data

Cooperative Multiagent Learning

Data Mining for Business

Mining Graphs, Trees and Sequences (MGTS 2005)

Knowledge Discovery and Ontologies (KDO 2005)

Knowledge Discovery from Data Streams

Reinforcement Learning in Non-stationary Environments

Discovery Challenge

Table of Contents

Invited Talks

Data Analysis in the Life Sciences — Sparking Ideas —	1
Machine Learning for Natural Language Processing (and Vice Versa?)	2
Statistical Relational Learning: An Inductive Logic Programming Perspective	3
Recent Advances in Mining Time Series Data	6
Focus the Mining Beacon: Lessons and Challenges from the World of E-Commerce	7
Data Streams and Data Synopses for Massive Data Sets	8

Long Papers

<i>k</i> -Anonymous Patterns	10
Interestingness is Not a Dichotomy: Introducing Softness in Constrained Pattern Mining	22
Generating Dynamic Higher-Order Markov Models in Web Usage Mining	34
TREE ² - Decision Trees for Tree Structured Data	46

Agglomerative Hierarchical Clustering with Constraints: Theoretical and Empirical Results	59
Cluster Aggregate Inequality and Multi-level Hierarchical Clustering	71
Ensembles of Balanced Nested Dichotomies for Multi-class Problems	84
Protein Sequence Pattern Mining with Constraints	96
An Adaptive Nearest Neighbor Classification Algorithm for Data Streams	108
Support Vector Random Fields for Spatial Classification	121
Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication	133
A Correspondence Between Maximal Complete Bipartite Subgraphs and Closed Patterns	146
Improving Generalization by Data Categorization	157
Mining Model Trees from Spatial Data	169
Word Sense Disambiguation for Exploiting Hierarchical Thesauri in Text Classification	181
Mining Paraphrases from Self-anchored Web Sentence Fragments	193
M ² SP: Mining Sequential Patterns Among Several Dimensions	205

A Systematic Comparison of Feature-Rich Probabilistic Classifiers for NER Tasks	217
Knowledge Discovery from User Preferences in Conversational Recommendation	228
Unsupervised Discretization Using Tree-Based Density Estimation	240
Weighted Average Pointwise Mutual Information for Feature Selection in Text Categorization	252
Non-stationary Environment Compensation Using Sequential EM Algorithm for Robust Speech Recognition	264
Hybrid Cost-Sensitive Decision Tree	274
Characterization of Novel HIV Drug Resistance Mutations Using Clustering, Multidimensional Scaling and SVM-Based Feature Ranking	285
Object Identification with Attribute-Mediated Dependences	297
Weka4WS: A WSRF-Enabled Weka Toolkit for Distributed Data Mining on Grids	309
Using Inductive Logic Programming for Predicting Protein-Protein Interactions from Multiple Genomic Data	321
ISOLLE: Locally Linear Embedding with Geodesic Distance	331

Active Sampling for Knowledge Discovery from Biomedical Data
..... 343

A Multi-metric Index for Euclidean and Periodic Matching
..... 355

Fast Burst Correlation of Financial Data
..... 368

A Propositional Approach to Textual Case Indexing
..... 380

A Quantitative Comparison of the Subgraph Miners MoFa, gSpan,
FFSM, and Gaston
..... 392

Efficient Classification from Multiple Heterogeneous Databases
..... 404

A Probabilistic Clustering-Projection Model for Discrete Data
..... 417

Short Papers

Collaborative Filtering on Data Streams
..... 429

The Relation of Closed Itemset Mining, Complete Pruning Strategies
and Item Ordering in Apriori-Based FIM Algorithms
..... 437

Community Mining from Multi-relational Networks
..... 445

Evaluating the Correlation Between Objective Rule Interestingness
Measures and Real Human Interest
..... 453

A Kernel Based Method for Discovering Market Segments in Beef Meat
..... 462

Corpus-Based Neural Network Method for Explaining Unknown Words by WordNet Senses	470
Segment and Combine Approach for Non-parametric Time-Series Classification	478
Producing Accurate Interpretable Clusters from High-Dimensional Data	486
Stress-Testing Hoeffding Trees	495
Rank Measures for Ordering	503
Dynamic Ensemble Re-Construction for Better Ranking	511
Frequency-Based Separation of Climate Signals	519
Efficient Processing of Ranked Queries with Sweeping Selection	527
Feature Extraction from Mass Spectra for Classification of Pathological States	536
Numbers in Multi-relational Data Mining	544
Testing Theories in Particle Physics Using Maximum Likelihood and Adaptive Bin Allocation	552
Improved Naive Bayes for Extremely Skewed Misclassification Costs	561
Clustering and Prediction of Mobile User Routes from Cellular Data	569

Elastic Partial Matching of Time Series
..... 577

An Entropy-Based Approach for Generating Multi-dimensional
Sequential Patterns 585

Visual Terrain Analysis of High-Dimensional Datasets
..... 593

An Auto-stopped Hierarchical Clustering Algorithm for Analyzing 3D
Model Database 601

A Comparison Between Block CEM and Two-Way CEM Algorithms to
Cluster a Contingency Table 609

An Imbalanced Data Rule Learner 617

Improvements in the Data Partitioning Approach for Frequent Itemsets
Mining 625

On-Line Adaptive Filtering of Web Pages 634

A Bi-clustering Framework for Categorical Data 643

Privacy-Preserving Collaborative Filtering on Vertically Partitioned
Data 651

Indexed Bit Map (IBM) for Mining Frequent Sequences 659

STochFS: A Framework for Combining Feature Selection Outcomes
Through a Stochastic Process 667

Speeding Up Logistic Model Tree Induction 675

A Random Method for Quantifying Changing Distributions in Data Streams	684
Deriving Class Association Rules Based on Levelwise Subspace Clustering	692
An Incremental Algorithm for Mining Generators Representation	701
Hybrid Technique for Artificial Neural Network Architecture and Weight Optimization	709
Author Index	717

Data Analysis in the Life Sciences

— Sparking Ideas —

Michael R. Berthold

ALTANA-Chair for Bioinformatics and Information Mining,
Dept. of Computer and Information Science, Konstanz University, Germany
`Michael.Berthold@uni-konstanz.de`

Data from various areas of Life Sciences have increasingly caught the attention of data mining and machine learning researchers. Not only is the amount of data available mind-boggling but the diverse and heterogenous nature of the information is far beyond any other data analysis problem so far. In sharp contrast to classical data analysis scenarios, the life science area poses challenges of a rather different nature for mainly two reasons. Firstly, the available data stems from heterogenous information sources of varying degrees of reliability and quality and is, without the interactive, constant interpretation of a domain expert, not useful. Furthermore, predictive models are of only marginal interest to those users – instead they hope for new insights into a complex, biological system that is only partially represented within that data anyway. In this scenario, the data serves mainly to create new insights and generate new ideas that can be tested. Secondly, the notion of feature space and the accompanying measures of similarity cannot be taken for granted. Similarity measures become context dependent and it is often the case that within one analysis task several different ways of describing the objects of interest or measuring similarity between them matter.

Some more recently published work in the data analysis area has started to address some of these issues. For example, data analysis in parallel universes [1], that is, the detection of patterns of interest in various different descriptor spaces at the same time, and mining of frequent, discriminative fragments in large, molecular data bases [2]. In both cases, sheer numerical performance is not the focus; it is rather the discovery of interpretable pieces of evidence that lights up new ideas in the users mind. Future work in data analysis in the life sciences needs to keep this in mind: the goal is to trigger new ideas and stimulate interesting associations.

References

1. Berthold, M.R., Wiswedel, B., Patterson, D.E.: Interactive exploration of fuzzy clusters using neighborgrams. *Fuzzy Sets and Systems* 149 (2005) 21–37
2. Hofer, H., Borgelt, C., Berthold, M.R.: Large scale mining of molecular fragments with wildcards. *Intelligent Data Analysis* 8 (2004) 376–385

Machine Learning for Natural Language Processing (and Vice Versa?)

Clai e Cardie

Department of Computer Science, Cornell University, USA
cardie@cs.cornell.edu
<http://www.cs.cornell.edu/home/cardie/>

Over the last 10-15 years, the influence of the field of machine learning has advanced the way that each is done in the field of natural language processing. This talk will begin by covering the history of this advancement. In particular, learning methods have proved successful in modeling word-and-alone extraction components to handle a number of linguistic tasks. Moreover, these components can be combined to model complex tasks that exhibit shallow extraction capabilities: they can, for example, extract key facts from unstructured documents in limited domains and answer general-domain questions from open-domain document collections. I will briefly describe the state of the art for these practical extraction applications, focusing on the role that machine learning methods have played in their development.

The second half of the talk will explore the role that natural language processing plays in machine learning research. Here, I will explain the kind of extracted features that are relatively easy to incorporate into machine learning datasets. In addition, I'll outline some notable foundations for natural language processing that have, in effect, opened the door for new machine learning algorithms.

Statistical Relational Learning: An Inductive Logic Programming Perspective

Lic De Raedt

Institute for Computer Science, Machine Learning Lab,
Albert-Ludwigs-University, Georges-Köhler-Allee, Gebäude 079,
D-79110 Freiburg i. Brg., Germany
deraedt@informatik.uni-freiburg.de

In the a few years, there have been a lot of working articles in the area of probabilistic logic programming, logic programming and machine learning [14,18,13,9,6,1,11]. This work is known under the name of statistical relational learning [7,5], probabilistic logic learning [4], or probabilistic inductive logic programming. Whereas most of the existing work have addressed formal probabilistic learning, expressive and extended probabilistic formalisms with relational aspects, I shall address a different expressive, in which I shall address formal inductive logic programming and study how inductive logic programming formalisms, learning and technique can be extended to deal with probabilistic inference. This addition has already combined a rich variety of valuable formalisms and techniques, including probabilistic Horn abduction by David Poole, PRISM by Sato, stochastic logic programming by Muggle [13] and Chen [2], Bayesian logic programming [10,8] by Keeling and De Raedt, and Logical Hidden Markov Model [11].

The main contribution of this article is the introduction of heterogeneous probabilistic inductive logic programming learning which are derived from the learning formalism, formal inference and formal proof learning of the field of inductive logic programming [3]. Each of the existing contributions differ in the notion of probabilistic logic inference, exact and probabilistic distribution. The learning, probabilistic learning formalism, is incorporated in the well-known PRISM by Sato [19] and Chen's 'Fail the Adjusted Maximization' approach of Sato et al. inference in stochastic logic programming [2]. A novel system has been recently developed and handles heterogeneous FOIL by Sato [12]. I combine the principles of the well-known inductive logic programming system FOIL [15] with the naive Bayesian approach. In probabilistic learning formalism, exact and ground facts should be probabilistically encoded by the logic programming. The second learning, probabilistic learning formalism, is incorporated in Bayesian logic programming [10,8], which in general Bayesian network with logic programming. This learning is also adopted by [6]. Exact learning are Hebb and inference should be a probabilistic model for the logic programming. The hybrid learning, learning formalism, is novel. It is motivated by the learning of stochastic context-free grammar for learning. In this learning, exact and ground facts should be probabilistically encoded for the well-known stochastic logic programming. The second learning (and hence inference) are by no means the only possible learning for probabilistic

inductive logic programming, built-in – I hope – provide a sufficient in-house infrastructure of his exciting field.

For a full survey of statistical relational learning on probabilistic inductive logic programming, see the forthcoming book of [4], and for a detailed account on the probabilistic inductive logic programming of [16], where a long and detailed evaluation of his contribution can be found.

Acknowledgements

This joint work with Kristian Kersting. The author would also like to thank Niel Landwehr and Sanna Toege for interesting collaboration on nFOIL and the learning of SLP, respectively. This work is a part of the EU IST FET/JCPC APRIL II (Application of Probabilistic Inductive Logic Programming II).

References

1. C. R. Anderson, P. Domingos, and D. S. Weld. Relational Markov Models and their Application to Adaptive Web Navigation. In D. Hand, D. Keim, O. R. Zaïne, and R. Goebel, editors, *Proceedings of the Eighth International Conference on Knowledge Discovery and Data Mining (KDD-02)*, pages 143–152, Edmonton, Canada, 2002. ACM Press.
2. J. Cussens. Loglinear models for first-order probabilistic reasoning. In K. B. Laskey and H. Prade, editors, *Proceedings of the Fifteenth Annual Conference on Uncertainty in Artificial Intelligence (UAI-99)*, pages 126–133, Stockholm, Sweden, 1999. Morgan Kaufmann.
3. L. De Raedt. Logical settings for concept-learning. *Artificial Intelligence*, 95(1):197–201, 1997.
4. L. De Raedt and K. Kersting. Probabilistic Logic Learning. *ACM-SIGKDD Explorations: Special issue on Multi-Relational Data Mining*, 5(1):31–48, 2003.
5. T. Dietterich, L. Getoor, and K. Murphy, editors. *Working Notes of the ICML-2004 Workshop on Statistical Relational Learning and its Connections to Other Fields (SRL-04)*, 2004.
6. N. Friedman, L. Getoor, D. Koller, and A. Pfeffer. Learning probabilistic relational models. In T. Dean, editor, *Proceedings of the Sixteenth International Joint Conferences on Artificial Intelligence (IJCAI-99)*, pages 1300–1309, Stockholm, Sweden, 1999. Morgan Kaufmann.
7. L. Getoor and D. Jensen, editors. *Working Notes of the IJCAI-2003 Workshop on Learning Statistical Models from Relational Data (SRL-03)*, 2003.
8. K. Kersting and L. De Raedt. Adaptive Bayesian Logic Programs. In C. Rouveirol and M. Sebag, editors, *Proceedings of the Eleventh Conference on Inductive Logic Programming (ILP-01)*, volume 2157 of *LNCS*, Strasbourg, France, 2001. Springer.
9. K. Kersting and L. De Raedt. Bayesian logic programs. Technical Report 151, University of Freiburg, Institute for Computer Science, April 2001.
10. K. Kersting and L. De Raedt. Towards Combining Inductive Logic Programming and Bayesian Networks. In C. Rouveirol and M. Sebag, editors, *Proceedings of the Eleventh Conference on Inductive Logic Programming (ILP-01)*, volume 2157 of *LNCS*, Strasbourg, France, 2001. Springer.

11. K. Kersting, T. Raiko, S. Kramer, and L. De Raedt. Towards discovering structural signatures of protein folds based on logical hidden markov models. In R. B. Altman, A. K. Dunker, L. Hunter, T. A. Jung, and T. E. Klein, editors, *Proceedings of the Pacific Symposium on Biocomputing*, pages 192 – 203, Kauai, Hawaii, USA, 2003. World Scientific.
12. N. Landwehr, K. Kersting, and L. De Raedt. nfoil: Integrating naive bayes and foil. In *Proceedings of the 20th National Conference on Artificial Intelligence*. AAAI Press, 2005.
13. S. H Muggleton. Stochastic logic programs. In L. De Raedt, editor, *Advances in Inductive Logic Programming*. IOS Press, 1996.
14. D. Poole. Probabilistic Horn abduction and Bayesian networks. *Artificial Intelligence*, 64:81–129, 1993.
15. J. R. Quinlan and R. M. Cameron-Jones. Induction of logic programs:FOIL and related systems. *New Generation Computing*, pages 287–312, 1995.
16. L. De Raedt and K. Kersting. Probabilistic inductive logic programming. In *Proceedings of the 15th International Conference on Algorithmic Learning Theory*. Springer, 2004.
17. L. De Raedt, K. Kersting, and S. Torge. Towards learning stochastic logic programs from proof-banks. In *Proceedings of the 20th National Conference on Artificial Intelligence*. AAAI Press, 2005.
18. T. Sato. A Statistical Learning Method for Logic Programs with Distribution Semantics. In L. Sterling, editor, *Proceedings of the Twelfth International Conference on Logic Programming (ICLP-1995)*, pages 715 – 729, Tokyo, Japan, 1995. MIT Press.
19. T. Sato and Y. Kameya. Parameter learning of logic programs for symbolic-statistical modeling. *Journal of Artificial Intelligence Research*, 15:391–454, 2001.

Recent Advances in Mining Time Series Data

Eamonn Keogh

Department of Computer Science & Engineering,
University of California, Riverside, USA

eamonn@cs.ucr.edu

<http://www.cs.ucr.edu/~eamonn>

Much of the world's supply of data is in the form of time series. For example, a web page, any type of data can be meaningfully converted into "time series", including text, DNA, video, images etc. The last decade has seen an explosion of interest in mining time series data for the academic community. There has been significant work on algorithms to classify, cluster, segment, index, discover, visualize, and detect anomalies/novelty in time series.

In this talk I will discuss the latest advance in mining time series data, including:

- New segmentation of time series data.
- New algorithm for denoising.
- The algorithm for a scalable online scalable...
- New analysis and application of time series data mining.

I will end the talk with a discussion of "what's left to do" in time series data mining.

References

1. E. Keogh. Exact indexing of dynamic time warping. In *Proceedings of the 8th International Conference on Very Large Data Bases*, pages 406–417, 2002.
2. E. Keogh and S. Kasetty. On the need for time series data mining benchmarks: A survey and empirical demonstration. In *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 102–111, 2002.
3. E. Keogh, J. Lin, and W. Truppel. Clustering of time series subsequences is meaningless: Implications for past and future research. In *Proceedings of the 3rd IEEE International Conference on Data Mining*, pages 115–122, 2003.
4. E. Keogh, S. Lonardi, and C. Ratanamahatana. Towards parameter-free data mining. In *Proceedings of the tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2004.
5. C.A. Ratanamahatana and E. Keogh. Everything you know about dynamic time warping is wrong. In *Proceedings of the Third Workshop on Mining Temporal and Sequential Data, in conjunction with the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD-2004)*, 2004.

Focus the Mining Beacon: Lessons and Challenges from the World of E-Commerce

Ron Kohavi

Microsoft Corporation, USA
ronnyk@cs.stanford.edu
<http://www.kohavi.com>

Electronic Commerce is now entering its second decade, with Amazon.com and eBay now in existence for ten years. With a diverse array of data, an actionable domain, and measurable ROI, little companies are data mining and knowledge discovery on demand and their success and innovation in action. We see an interesting lesson and challenge in using e-commerce example applications: (i) business-level technical, and (ii) the mining lifecycle for data collection, data warehousing, consolidation, discovery and deployment. Many of the lessons and challenges are applicable to domain outside e-commerce.

Data Streams and Data Synopses for Massive Data Sets (Invited Talk)

Yossi Matias

Tel Aviv University,
HyperRoll Inc., Stanford University
matias@tau.ac.il

Abstract. With the proliferation of data intensive applications, it has become necessary to develop new techniques to handle massive data sets. Traditional algorithmic techniques and data structures are not always suitable to handle the amount of data that is required and the fact that the data often streams by and cannot be accessed again. A field of research established over the past decade is that of handling massive data sets using data synopses, and developing algorithmic techniques for data stream models. We will discuss some of the research work that has been done in the field, and provide a decades' perspective to data synopses and data streams.

1 Summary

In recent years, we have witnessed an explosion in data used in various applications. In general, the growth rate in data is known to exceed the increase rate in the size of RAM, and of the available computation power (a.k.a. Moore's Law). As a result, traditional algorithms and data structures are often no longer adequate to handle the massive data sets required by these applications.

One approach to handle massive data sets is to use *external memory algorithms*, designed to make an effective utilization of I/O. In such algorithms the data structures are often implemented in external storage devices, and the objective is in general to minimize the number of I/Os. For a survey of works on external memory algorithms see [6]. Such algorithms assume that the entire input data is available for further processing. There are, however, many applications where the data is only seen once, as it "streams by". This may be the case in, e.g., financial applications, network monitoring, security, telecommunications data management, web applications, manufacturing, and sensor networks. Even in data warehouse applications, where the data may in general be available for additional querying, there are many situations where data analysis needs to be done as the data is loaded into the data warehouse, since the cost of accessing the data in a fully loaded production system may be significantly larger than just the basic cost of I/O. Additionally, even in the largest data warehouses, consisting of hundreds of terabytes, data is only maintained for a limited time, so access to historical data may often be infeasible.

It had thus become necessary to address situations in which massive data sets are required to be handled as they "stream by", and using only limited memory. Motivated by this need, the research field of data streams and data synopses has emerged

and established over the last few years. We will discuss some of the research work that has been done in the field, and provide a decades' perspective to data streams and data synopses. A longer version of this abstract will be available at [4].

The data stream model is quite simple: it is assumed that the input data set is given as a sequence of data items. Each data item is seen only once, and any computation can be done utilizing the data structures maintained in main memory. These memory resident data structures are substantially smaller than the input data. As such, they cannot fully represent the data as is the case for traditional data structures, but can only provide a synopsis of the input data; hence they are denoted as *synopsis data structures*, or *data synopses* [3].

The use of data synopses implies that data analysis that is dependent on the entire streaming data will often be approximated. Furthermore, ad hoc queries that are dependent on the entire input data could only be served by the data synopses, and as a result only approximate answers to queries will be available. A primary objective in the design of data synopses is to have the smallest data synopses that would guarantee small, and if possible bounded, error on the approximated computation.

As we have shown in [1], some essential statistical data analysis, the so-called *frequency moments*, can be approximated using synopses that are as small as polynomial or even logarithmic in the input size. Over the last few years there has been a proliferation of additional works on data streams and data synopses. See, e.g., the surveys [2] and [5]. These works include theoretical results, as well as applications in databases, network traffic analysis, security, sensor networks, and program profiling; synopses include samples, random projections, histograms, wavelets, and XML synopses, among others. There remain a plethora of interesting open problems, both theoretical as well as applied.

References

1. Alon, N., Matias, Y., Szegedy, M.: The space complexity of approximating the frequency moments. *J. of Computer and System Sciences* 58 (1999), 137-147. STOC'96 Special Issue
2. Babcock, B., Babu, S., Datar, M., Motwani, R., Widom, J.: Models and issues in data stream systems. In *Proc. Symposium on Principles of Database Systems* (2002), 1-16
3. Gibbons, P.B., Matias, Y.: Synopses data structures for massive data sets. In: *External memory algorithms*, DIMACS Series Discrete Math. & TCS, AMS, 50 (1999). Also SODA'99
4. Matias, Y.: Data streams and data synopses for massive data sets.
<http://www.cs.tau.ac.il/~matias/streams/>
5. Muthukrishnan, S.: Data streams: Algorithms and applications.
<http://www.cs.rutgers.edu/~muthu/stream-1-1.ps>
6. Vitter, J.S.: External memory algorithms and data structures. *ACM Comput Surv.* 33(2): 209-271 (2001)

k-Anonymous Patterns

Maurizio Atzori^{1,2}, Francesco Bonchi², Fosca Giannotti², and Dino Pedreschi¹

¹ Pisa KDD Laboratory, Computer Science Department, University of Pisa, Italy
{atzori, pedre}@di.unipi.it

² Pisa KDD Laboratory, ISTI - CNR, Pisa, Italy
{francesco.bonchi, fosca.giannotti}@isti.cnr.it

Abstract. It is generally believed that data mining results do not violate the *anonymity* of the individuals recorded in the source database. In fact, data mining models and patterns, in order to ensure a required statistical significance, represent a large number of individuals and thus conceal individual identities: this is the case of the minimum support threshold in association rule mining. In this paper we show that this belief is ill-founded. By shifting the concept of *k-anonymity* from data to patterns, we formally characterize the notion of a threat to anonymity in the context of pattern discovery, and provide a methodology to efficiently and effectively identify all possible such threats that might arise from the disclosure of a set of extracted patterns.

1 Introduction

Privacy Preserving Data Mining, i.e., the analysis of data mining side-effects on privacy, has recently become a key research issue and is receiving a growing attention from the research community [1,3,9,16]. However, despite such efforts, a common understanding of what is meant by “privacy” is still missing. This fact has led to the proliferation of many completely different approaches to privacy preserving data mining, all sharing the same generic goal: producing a valid mining model without disclosing “private” data. As highlighted in [9], the approaches pursued so far leave a privacy question open: do the data mining results themselves violate privacy? Put in other words, do the disclosure of extracted patterns open up the risk of privacy breaches that may reveal sensitive information? During the last year, few works [7,9,11] have tried to address this problem by some different points of view, but they all require some *a priori* knowledge of what is sensitive and what is not.

In this paper we study when data mining results represent *per se* a threat to privacy, independently of any background knowledge of what is sensitive. In particular, we focus on *individual privacy*, which is concerned with the *anonymity* of individuals.

A prototypical application instance is in the medical domain, where the collected data are typically very sensitive, and the kind of privacy usually required is the anonymity of the patients in a survey. Consider a medical institution where the usual hospital activity is coupled with medical research activity. Since physicians are the data collectors and holders, and they already know everything about their patients, they have unrestricted access to the collected information. Therefore, they can perform real mining on all available information using traditional mining tools – not necessarily the privacy preserving

ones. This way they maximize the outcome of the knowledge discovery process, without any concern about privacy of the patients which are recorded in the data. But the anonymity of patients becomes a key issue when the physicians want to share their discoveries (e.g., association rules holding in the data) with their scientific community.

At a first sight, it seems that data mining results do not violate the anonymity of the individuals recorded in the source database. In fact, data mining models and patterns, in order to ensure a required statistical significance, represent a large number of individuals and thus conceal individual identities: this is the case of the minimum support threshold in association rule mining. The next example shows that the above belief is ill-founded.

Example 1. Consider the following association rule:

$$a_1 \wedge a_2 \wedge a_3 \Rightarrow a_4 \quad [sup = 80, conf = 98.7\%]$$

where *sup* and *conf* are the usual interestingness measures of *support* and *confidence* as defined in [2]. Since the given rule holds for a number of individuals (80), which seems large enough to protect individual privacy, one could conclude that the given rule can be safely disclosed. But, is this all the information contained in such a rule? Indeed, one can easily derive the support of the premise of the rule:

$$sup(\{a_1, a_2, a_3\}) = \frac{sup(\{a_1, a_2, a_3, a_4\})}{conf} \approx \frac{80}{0.987} = 81.05$$

Given that the pattern $a_1 \wedge a_2 \wedge a_3 \wedge a_4$ holds for 80 individuals, and that the pattern $a_1 \wedge a_2 \wedge a_3$ holds for 81 individuals, we can infer that in our database there is just one individual for which the pattern $a_1 \wedge a_2 \wedge a_3 \wedge \neg a_4$ holds.

The knowledge inferred is a clear threat to the anonymity of that individual: on one hand the pattern identifying the individual could itself contain sensitive information; on the other hand it could be used to re-identify the same individual in other databases.

It is worth noting that this problem is very general: the given rule could be, instead of an association, a classification rule, or the path from the root to the leaf in a decision tree, and the same reasoning would still hold. Moreover, it is straightforward to note that, unluckily, the more accurate is a rule, the more unsafe it may be w.r.t. anonymity. As shown later, this anonymity problem can not be simply solved by discarding the most accurate rules: in fact, more complex kinds of threats to anonymity exist which involve more than simply two itemsets.

1.1 Related Works

During the last years a novel problem has emerged in privacy-preserving data mining [7,9,11]: do the data mining results themselves violate privacy? Only little preliminary work is available. The work in [9] studies the case of a classifier trained over a mixture of different kind of data: *public* (known to every one including the adversary), *private/sensitive* (should remain unknown to the adversary), and *unknown* (neither sensitive nor known by the adversary). The authors propose a model for privacy implication of the learned classifier.

In [11] the data owner, rather than sharing the data, prefers to share the mined association rules, but requires that a set of *restricted* association rules are not disclosed. The authors propose a framework to sanitize the output of association rules mining, while blocking some inference channels for the restricted rules.

In [7] a framework for evaluating classification rules in terms of their perceived privacy and ethical sensitivity is described. The proposed framework empowers the data miner with alerts for sensitive rules that can be accepted or dismissed by the user as appropriate. Such alerts are based on an aggregate *sensitivity combination function*, which assigns to each rule a value of sensitivity by aggregating the sensitivity value (an integer between 0 and 9) of each attribute involved in the rule. The process of labelling each attribute with its sensitivity value must be accomplished by the domain expert.

The fundamental difference of these approaches with ours lies in generality: we propose a novel, objective definition of privacy compliance of patterns without any reference to a preconceived knowledge of sensitive data or patterns, on the basis of the rather intuitive and realistic constraint that the anonymity of individuals should be guaranteed.

An important method for protecting individual privacy is *k-anonymity*, introduced in [14], a notion that establishes that the cardinality of the answer to any possible query should be at least k . In this work, it is shown that protection of individual sources does not guarantee protection when sources are cross-examined: a sensitive medical record, for instance, can be uniquely linked to a *named* voter record in a publicly available voter list through some shared attributes. The objective of *k-anonymity* is to eliminate such opportunities of inferring private information through cross linkage. In particular, this is obtained by a “sanitization” of the source data that is transformed in such a way that, for all possible queries, at least k tuples will be returned. Such a sanitization is obtained by generalization and suppression of attributes and/or tuples [15].

Trivially, by mining a *k-anonymized* database no patterns threatening the anonymity can be obtained. But such mining would produce models impoverished by the information loss which is intrinsic in the generalization and suppression techniques. Since our objective is to extract valid and interesting patterns, we propose to postpone *k-anonymization* after the actual mining step. In other words, we do not to enforce *k-anonymity* onto the source data, but instead we move such a concept to the extracted patterns.

1.2 Paper Contributions

In this paper we study the privacy problem described above in the very general setting of patterns which are boolean formulas over a binary database. Our contribution is twofold:

- we define *k-anonymous* patterns and provide a general characterization of inference channels holding among patterns that may threaten anonymity of source data;
- we develop an effective and efficient algorithm to detect such potential threats, which yields a methodology to check whether the mining results may be disclosed without any risk of violating anonymity.

We emphasize that the capability of detecting the potential threats is extremely useful for the analyst to determine a trade-off among the quality of mining result and the

privacy guarantee, by means of an iterative interaction with the proposed detection algorithm. Our empirical experiments, reported in this paper, bring evidence to this observation. It should also be noted the different setting w.r.t. the other works in privacy preserving data mining: in our context no data perturbation or sanitization is performed; we allow real mining on the real data, while focussing on the *anonymity preservation properties of the extracted patterns*. We have also developed possible strategies to eliminate the threats to anonymity by introducing distortion on the dangerous patterns in a controlled way: for lack of space these results are omitted here but can be found in [5].

2 k -Anonymous Patterns and σ -Frequent Itemsets

We start by defining binary databases and patterns following the notation in [8].

Definition 1. A binary database $\mathcal{D} = (\mathcal{I}, \mathcal{T})$ consists of a finite set of binary variables $\mathcal{I} = \{i_1, \dots, i_p\}$, also known as items, and a finite multiset $\mathcal{T} = \{t_1, \dots, t_n\}$ of p -dimensional binary vectors recording the values of the items. Such vectors are also known as transactions. A pattern for the variables in \mathcal{I} is a logical (propositional) sentence built by AND (\wedge), OR (\vee) and NOT (\neg) logical connectives, on variables in \mathcal{I} . The domain of all possible patterns is denoted $\mathcal{Pat}(\mathcal{I})$.

According to Def. 1, $e \wedge (\neg b \vee \neg d)$, where $b, d, e \in \mathcal{I}$, is a pattern. One of the most important properties of a pattern is its frequency in the database, i.e. the number of individuals (transactions) in the database which make the given pattern true¹.

Definition 2. Given a database \mathcal{D} , a transaction $t \in \mathcal{D}$ and a pattern p , we write $p(t)$ if t makes p true. The support of p in \mathcal{D} is given by the number of transactions which makes p true: $\text{supp}_{\mathcal{D}}(p) = |\{t \in \mathcal{D} \mid p(t)\}|$.

The most studied *pattern class* is the itemset, i.e., a conjunction of positive valued variables, or in other words, a set of items. The retrieval of itemsets which satisfy a minimum frequency property is the basic step of many data mining tasks, including (but not limited to) association rules [2,4].

Definition 3. The set of all itemsets $2^{\mathcal{I}}$, is a pattern class consisting of all possible conjunctions of the form $i_1 \wedge i_2 \wedge \dots \wedge i_m$. Given a database \mathcal{D} and a minimum support threshold σ , the set of σ -frequent itemsets in \mathcal{D} is denoted

$$\mathcal{F}(\mathcal{D}, \sigma) = \{ \langle X, \text{supp}_{\mathcal{D}}(X) \rangle \mid X \in 2^{\mathcal{I}} \wedge \text{supp}_{\mathcal{D}}(X) \geq \sigma \}$$

Itemsets are usually denoted in the form of set of the items in the conjunction, e.g. $\{i_1, \dots, i_m\}$; or sometimes, simply $i_1 \dots i_m$. Figure 1(b) shows the different notation used for general patterns and for itemsets. The problem addressed in this paper is given by the possibility of inferring from the output of frequent itemset mining, i.e. $\mathcal{F}(\mathcal{D}, \sigma)$, the existence of patterns with very low support (i.e., smaller than an anonymity threshold k , but not null): such patterns represent a threat for the anonymity of the individuals about which they are true.

¹ The notion of truth of a pattern w.r.t. a transaction t is defined in the usual way: t makes p true iff t is a model of the propositional sentence p .

\mathcal{D}		<p>Notation: patterns</p> $sup_{\mathcal{D}}(a \vee f) = 11$ $sup_{\mathcal{D}}(e \wedge (\neg b \vee \neg d)) = 4$ $sup_{\mathcal{D}}(h \wedge \neg b) = 1$ <p>Notation: itemsets</p> $sup_{\mathcal{D}}(abc) = 6$ $sup_{\mathcal{D}}(abde) = 7$ $sup_{\mathcal{D}}(cd) = 9$	$\mathcal{F}(\mathcal{D}, 8) = \{\langle \emptyset, 12 \rangle, \langle a, 9 \rangle, \langle b, 8 \rangle, \langle c, 9 \rangle, \langle d, 10 \rangle, \langle e, 11 \rangle, \langle ab, 8 \rangle, \langle ae, 8 \rangle, \langle cd, 9 \rangle, \langle ce, 9 \rangle, \langle de, 10 \rangle, \langle cde, 9 \rangle\}$ <p style="text-align: center;">(c)</p>	
(a)	(b)			(d)
				(e)
t_1	1 1 1 1 1 1 1 1 1			$Cl(\mathcal{D}, 8) = \{\langle \emptyset, 12 \rangle, \langle a, 9 \rangle, \langle e, 11 \rangle, \langle ab, 8 \rangle, \langle ae, 8 \rangle, \langle de, 10 \rangle, \langle cde, 9 \rangle\}$ <p style="text-align: center;">(d)</p>
t_2	1 1 1 1 1 0 1 0			
t_3	1 1 1 1 1 0 0 0			
t_4	1 1 1 1 1 1 1 0			
t_5	1 1 1 1 1 0 0 0			
t_6	1 1 1 1 1 0 0 0			
t_7	1 1 0 1 1 0 0 0			
t_8	1 0 0 0 1 1 1 0			
t_9	0 0 1 1 1 1 1 0			
t_{10}	0 0 1 1 1 0 0 0			
t_{11}	0 0 1 1 1 1 1 1			
t_{12}	1 1 0 0 0 1 1 0			

Fig. 1. Running example: (a) the binary database \mathcal{D} ; (b) different notation used for patterns and itemsets; (c) the set of σ -frequent ($\sigma = 8$) itemsets; (d) the set of closed frequent itemsets; (e) the set of maximal inference channels for $k = 3$ and $\sigma = 6$

Definition 4. Given a database \mathcal{D} and an anonymity threshold k , a pattern p is said to be k -anonymous if $sup_{\mathcal{D}}(p) \geq k$ or $sup_{\mathcal{D}}(p) = 0$.

2.1 Problem Definition

Before introducing our anonymity preservation problem, we need to define the inference of supports, which is the basic tool for the attacks to anonymity.

Definition 5. A set S of pairs $\langle X, n \rangle$, where $X \in 2^{\mathcal{I}}$ and $n \in \mathbb{N}$, and a database \mathcal{D} are said to be σ -compatible if $S \subseteq \mathcal{F}(\mathcal{D}, \sigma)$. Given a pattern p we say that $S \models_{\mathcal{D}}(p) > x$ (respectively $S \models_{\mathcal{D}}(p) < x$) if, for all databases \mathcal{D} σ -compatible with S , we have that $sup_{\mathcal{D}}(p) > x$ (respectively $sup_{\mathcal{D}}(p) < x$).

Informally, we call *inference channel* any subset of the collection of itemsets (with their respective supports), from which it is possible to infer non k -anonymous patterns. Our mining problem can be seen as frequent pattern extraction with two frequency thresholds: the usual minimum support threshold σ for itemsets (as defined in Definition 3), and an anonymity threshold k for general patterns (as defined in Definition 1).

Note that an itemset with support less than k is itself a non k -anonymous, and thus dangerous, pattern. However, since we can safely assume (as we will do in the rest of this paper) that $\sigma \gg k$, such pattern would be discarded by the usual mining algorithms.

Definition 6. Given a collection of frequent itemsets $\mathcal{F}(\mathcal{D}, \sigma)$ and an anonymity threshold k , our problem consists in detecting all possible inference channels $\mathcal{C} \subseteq \mathcal{F}(\mathcal{D}, \sigma) : \exists p \in \mathcal{P}at(\mathcal{I}) : \mathcal{C} \models 0 < sup_{\mathcal{D}}(p) < k$.

Obviously, a solution to this problem directly yields a method to formally prove that the disclosure of a given collection of frequent itemsets does not violate the anonymity

constraint: it is sufficient to check that no inference channel exists for the given collection. In this case, the collection can be safely distributed even to malicious adversaries. On the contrary, if this is not the case, we can proceed in two ways:

- mine a new collection of frequent itemsets under different circumstances, e.g., higher minimum support threshold, to look for an admissible collection;
- transform (sanitize) the collection to remove the inference channels.

The second alternative opens up many interesting mining problems, which are omitted here for lack of space, and are discussed in [5].

3 Detecting Inference Channels

In this Section we study how information about non k -anonymous patterns can be possibly inferred from a collection of σ -frequent itemsets. As suggested by Example 1, a simple inference channel is given by any itemset X which has a superset $X \cup \{a\}$ such that $0 < \sigma_{\mathcal{D}}(X) - \sigma_{\mathcal{D}}(X \cup \{a\}) < k$. In this case the pair $\langle X, \sigma_{\mathcal{D}}(X) \rangle, \langle X \cup \{a\}, \sigma_{\mathcal{D}}(X \cup \{a\}) \rangle$ is an inference channel for the non k -anonymous pattern $X \wedge \neg a$, whose support is directly given by $\sigma_{\mathcal{D}}(X) - \sigma_{\mathcal{D}}(X \cup \{a\})$. This is a trivial kind of inference channel. Do more complex structures of itemsets exist that can be used as inference channels? In general, the support of a pattern $p = i_1 \wedge \dots \wedge i_m \wedge \neg a_1 \wedge \dots \wedge \neg a_n$ can be inferred if we know the support of itemsets $I = \{i_1, \dots, i_m\}$, $J = I \cup \{a_1, \dots, a_n\}$, and every itemset L such that $I \subset L \subset J$.

Lemma 1. *Given a pattern $p = i_1 \wedge \dots \wedge i_m \wedge \neg a_1 \wedge \dots \wedge \neg a_n$ we have that:*

$$\sigma_{\mathcal{D}}(p) = \sum_{I \subset X \subset J} (-1)^{|X \setminus I|} \sigma_{\mathcal{D}}(X)$$

where $I = \{i_1, \dots, i_m\}$ and $J = I \cup \{a_1, \dots, a_n\}$.

Proof. (Sketch) The proof follows directly from the definition of support and the well-known *inclusion-exclusion principle* [10].

Following the notation in [6], we denote the right-hand side of the equation above as $f_I^J(\mathcal{D})$. In the database \mathcal{D} in Figure 1 we have that $\sigma_{\mathcal{D}}(b \wedge \neg d \wedge \neg e) = f_b^{bde}(\mathcal{D}) = \sigma_{\mathcal{D}}(b) - \sigma_{\mathcal{D}}(bd) - \sigma_{\mathcal{D}}(be) + \sigma_{\mathcal{D}}(bde) = 8 - 7 - 7 + 7 = 1$.

Definition 7. *Given a database \mathcal{D} , and two itemsets $I, J \in 2^{\mathcal{I}}$, $I = \{i_1, \dots, i_m\}$ and $J = I \cup \{a_1, \dots, a_n\}$, if $0 < f_I^J(\mathcal{D}) < k$, then the set of itemsets $\{X \mid I \subseteq X \subseteq J\}$ constitutes an inference channel for the non k -anonymous pattern $p = i_1 \wedge \dots \wedge i_m \wedge \neg a_1 \wedge \dots \wedge \neg a_n$. We denote such inference channel \mathcal{C}_I^J and we write $\sigma_{\mathcal{D}}(\mathcal{C}_I^J) = f_I^J(\mathcal{D})$.*

Example 2. Consider the database \mathcal{D} of Figure 1, and suppose $k = 3$. We have that \mathcal{C}_{ab}^{abcde} is an inference channel of support 1. This means that there is only one transaction $t \in \mathcal{D}$ is such that $a \wedge b \wedge \neg c \wedge \neg d \wedge \neg e$.

The next Theorem states that if there exists a non k -anonymous pattern, then there exists a pair of itemsets $I \subseteq J \in 2^{\mathcal{I}}$ such that \mathcal{C}_I^J is an inference channel.

Theorem 1. $\forall p \in \mathcal{Pat}(\mathcal{I}) : 0 < \dots, \mathcal{D}(p) < k . \exists I \subseteq J \in 2^{\mathcal{I}} : C_I^J$.

Proof. The case of a conjunctive pattern p is a direct consequence of Lemma 1. Let us now consider a generic pattern $p \in \mathcal{Pat}(\mathcal{I})$. Without loss of generality p is in *normal disjunctive form*: $p = p_1 \vee \dots \vee p_q$, where $p_1 \dots p_q$ are conjunctive patterns. We have that:

$$\dots, \mathcal{D}(p) \geq \max_{1 \leq i \leq q} \dots, \mathcal{D}(p_i).$$

Since $\dots, \mathcal{D}(p) < k$ we have for all patterns p_i that $\dots, \mathcal{D}(p_i) < k$. Moreover, since $\dots, \mathcal{D}(p) > 0$ is there at least a pattern p_i such that $\dots, \mathcal{D}(p_i) > 0$. Therefore, there is at least a conjunctive pattern p_i such that $0 < \dots, \mathcal{D}(p_i) < k$.

From Theorem 1 we conclude that all possible threats to anonymity are due to inference channels of the form C_I^J . However we can divide such inference channels in two subgroups:

1. inference channels involving only frequent itemsets;
2. inference channels involving also infrequent itemsets.

The first problem, addressed in the rest of this paper, is the most essential. In fact, a malicious adversary can easily find inference channels made up only of elements which are present in the disclosed output. However, these inference channels are not the unique possible source of inference: further inference channels involving also infrequent itemsets could be possibly discovered, albeit in a much more complex way.

In fact, in [6] deduction rules to derive tight bounds on the support of itemsets are introduced. Given an itemset J , if for each subset $I \subset J$ the support $\dots, \mathcal{D}(I)$ is known, such rules allow to compute lower and upper bounds on the support of J . Let l be the greatest lower bound we can derive, and u the smallest upper bound we can derive: if we find that $l = u$ then we can infer that $\dots, \mathcal{D}(J) = l = u$ without actual counting. In this case J is said to be a *derivable itemset*. We transpose such deduction techniques in our context and observe that they can be exploited to discover information about infrequent itemsets, and from these to infer non k -anonymous patterns. For lack of space, this higher-order problem is not discussed here, and left to the extended version of this paper. However, here we can say that the techniques to detect this kind of inference channels and to block them are very similar to the techniques for the first kind of channels. This is due to the fact that both kinds of channels rely on the same concept: inferring supports of larger itemsets from smaller ones. Indeed, the key equation of our work (Lemma 1) is also the basis of the deduction rules proposed in [6].

From now on we restrict our attention to the essential form of inference channel, namely those involving frequent itemsets only.

Definition 8. *The set of all C_I^J holding in $\mathcal{F}(\mathcal{D}, \sigma)$, together with their supports, is denoted $Ch(k, \mathcal{F}(\mathcal{D}, \sigma)) = \{ \langle C_I^J, f_I^J(\mathcal{D}) \rangle \mid 0 < f_I^J(\mathcal{D}) < k \wedge \langle J, \dots, \mathcal{D}(J) \rangle \in \mathcal{F}(\mathcal{D}, \sigma) \}$.*

Algorithm 1 detects all possible inference channels $Ch(k, \mathcal{F}(\mathcal{D}, \sigma))$ that hold in a collection of frequent itemsets $\mathcal{F}(\mathcal{D}, \sigma)$ by checking all possible pairs of itemsets $I, J \in \mathcal{F}(\mathcal{D}, \sigma)$ such that $I \subseteq J$. This could result in a very large number of checks. Suppose that $\mathcal{F}(\mathcal{D}, \sigma)$ is formed only by a maximal itemset Y and all its subsets (an

Algorithm 1 Naïve Inference Channel Detector

Input: $\mathcal{F}(\mathcal{D}, \sigma), k$
Output: $Ch(k, \mathcal{F}(\mathcal{D}, \sigma))$

- 1: $Ch(k, \mathcal{F}(\mathcal{D}, \sigma)) = \emptyset$
- 2: **for all** $\langle J, sup(J) \rangle \in \mathcal{F}(\mathcal{D}, \sigma)$ **do**
- 3: **for all** $I \subseteq J$ **do**
- 4: **compute** f_I^J ;
- 5: **if** $0 < f_I^J < k$ **then**
- 6: **insert** $\langle \mathcal{C}_I^J, f_I^J \rangle$ **in** $Ch(k, \mathcal{F}(\mathcal{D}, \sigma))$;

itemset is maximal if none of its proper supersets is in $\mathcal{F}(\mathcal{D}, \sigma)$). If $|Y| = n$ we get $|\mathcal{F}(\mathcal{D}, \sigma)| = 2^n$ (we also count the empty set), while the number of possible \mathcal{C}_I^J is $\sum_{1 \leq i \leq n} \binom{n}{i} (2^i - 1)$. In the following Section we study some interesting properties that allow to dramatically reduce the number of checks needed to retrieve $Ch(k, \mathcal{F}(\mathcal{D}, \sigma))$.

4 A Condensed Representation of Inference Channels

In this section we introduce a condensed representation of $Ch(k, \mathcal{F}(\mathcal{D}, \sigma))$, i.e., a subset of $Ch(k, \mathcal{F}(\mathcal{D}, \sigma))$ which is more efficient to compute, and sufficient to reconstruct the whole $Ch(k, \mathcal{F}(\mathcal{D}, \sigma))$. The benefits of having such condensed representation go far beyond mere efficiency. In fact, removing the redundancy existing in $Ch(k, \mathcal{F}(\mathcal{D}, \sigma))$, we also implicitly avoid redundant sanitization, when we will block inference channels holding in $\mathcal{F}(\mathcal{D}, \sigma)$ (recall that, as stated before, the issue of how to block inference channels is not covered in this paper).

Consider, for instance, the two inference channels $\langle \mathcal{C}_{ad}^{acd}, 1 \rangle$ and $\langle \mathcal{C}_{abd}^{abcd}, 1 \rangle$ holding in the database in Fig. 1(a): one is more specific than the other, but they both uniquely identify transaction t_7 . It is easy to see that many other families of equivalent, and thus redundant, inference channels can be found. *How can we directly identify one and only one representative inference channel in each family of equivalent ones?* The theory of *closed itemsets* can help us with this problem.

Closed itemsets were first introduced in [12] and since then they have received a great deal of attention especially by an algorithmic point of view [17,13]. They are a concise and lossless representation of all frequent itemsets, i.e., they contain the same information without redundancy. Intuitively, a closed itemset groups together all its subsets that have its same support; or in other words, it groups together itemsets which identify the same group of transactions.

Definition 9. Given the function $f(T) = \{i \in \mathcal{I} \mid \forall t \in T, i \in t\}$, which returns all the items included in the set of transactions T , and the function $g(X) = \{t \in \mathcal{T} \mid \forall i \in X, i \in t\}$ which returns the set of transactions supporting a given itemset X , the composite function $c = f \circ g$ is the closure operator. An itemset I is closed iff and only if $c(I) = I$. Given a database \mathcal{D} and a minimum support threshold σ , the set of frequent closed itemsets is denoted $Cl(\mathcal{D}, \sigma)$. An itemset $I \in Cl(\mathcal{D}, \sigma)$ is said to be maximal iff $\nexists J \supset I$ s.t. $J \in Cl(\mathcal{D}, \sigma)$.

Analogously to what happens for the pattern class of itemsets, if we consider the pattern class of conjunctive patterns we can rely on the *anti-monotonicity property of frequency*. For instance, the number of transactions for which the pattern $a \wedge \neg c$ holds is always larger than the number of transactions for which the pattern $a \wedge b \wedge \neg c \wedge \neg d$ holds.

Definition 10. Given two inference channels \mathcal{C}_I^J and \mathcal{C}_H^L we say that $\mathcal{C}_I^J \preceq \mathcal{C}_H^L$ when $I \subseteq H$ and $(J \setminus I) \subseteq (L \setminus H)$.

Proposition 1. $\mathcal{C}_I^J \preceq \mathcal{C}_H^L \Rightarrow \forall \mathcal{D} . f_I^J(\mathcal{D}) \geq f_H^L(\mathcal{D})$.

Therefore, when detecting inference channels, whenever we find a \mathcal{C}_H^L such that $f_H^L(\mathcal{D}) \geq k$, we can avoid checking the support of all inference channels $\mathcal{C}_I^J \preceq \mathcal{C}_H^L$, since they will be k -anonymous.

Definition 11. An inference channel \mathcal{C}_I^J is said to be maximal w.r.t. \mathcal{D} and σ , if $\forall H, L$ such that $I \subseteq H$ and $(J \setminus I) \subseteq (L \setminus H)$, $f_H^L = 0$. The set of maximal inference channels is denoted $\mathcal{MCh}(k, Cl(\mathcal{D}, \sigma))$.

Proposition 2. $\mathcal{C}_I^J \in \mathcal{MCh}(k, Cl(\mathcal{D}, \sigma)) \Rightarrow I \in Cl(\mathcal{D}, \sigma) \wedge J$ is maximal.

Proof. *i)* $I \in Cl(\mathcal{D}, \sigma)$: if I is not closed then consider its closure $c(I)$ and consider $J' = J \cup (c(I) \setminus I)$. For the definition of closure, the set of transactions containing I is the same of the set of transactions containing $c(I)$, and the set of transactions containing J' is the same of the set of transactions containing J . It follows that $\mathcal{C}_{c(I)}^{J'} \succeq \mathcal{C}_I^J$ and $f_{c(I)}^{J'} = f_I^J > 0$. Then, if I is not closed, \mathcal{C}_I^J is not maximal.

ii) J is maximal: if J is not maximal then consider its frequent superset $J' = J \cup \{a\}$ and consider $I' = I \cup a$. It is straightforward to see that $f_I^J = f_I^{J'} + f_{I'}^{J'}$ and that $\mathcal{C}_I^{J'} \succeq \mathcal{C}_I^J$ and $\mathcal{C}_{I'}^{J'} \succeq \mathcal{C}_I^J$. Therefore, since $f_I^J > 0$, at least one among $f_{I'}^{J'}$ and $f_I^{J'}$ must be not null. Then, if J is not maximal, \mathcal{C}_I^J is not maximal as well.

The next Theorem shows how the support of any channel in $\mathcal{Ch}(k, \mathcal{F}(\mathcal{D}, \sigma))$ can be reconstructed from $\mathcal{MCh}(k, Cl(\mathcal{D}, \sigma))$.

Theorem 2. Given $\mathcal{C}_I^J \in \mathcal{Ch}(k, \mathcal{F}(\mathcal{D}, \sigma))$, let M be any maximal itemset such that $M \supseteq J$. The following equation holds:

$$f_I^J(\mathcal{D}) = \sum_{c(X)} f_{c(X)}^M(\mathcal{D})$$

where $c(I) \subseteq c(X) \subseteq M$ and $c(X) \cap (J \setminus I) = \emptyset$.

Proof. See [5].

From Theorem 2 we conclude that all the addends needed to compute $f_I^J(\mathcal{D})$ for an inference channel are either in $\mathcal{MCh}(k, Cl(\mathcal{D}, \sigma))$ or are null. Therefore, as the set of all closed frequent itemsets $Cl(\mathcal{D}, \sigma)$ contains all the information of $\mathcal{F}(\mathcal{D}, \sigma)$ in a more compact representation, analogously the set $\mathcal{MCh}(k, Cl(\mathcal{D}, \sigma))$ represents, without redundancy, all the information in $\mathcal{Ch}(k, \mathcal{F}(\mathcal{D}, \sigma))$.

In the database \mathcal{D} of Figure 1(a), given $\sigma = 6$ and $k = 3$, $|Ch(3, \mathcal{F}(\mathcal{D}, 6))| = 58$ while $|MCh(3, Cl(\mathcal{D}, 6))| = 5$ (Figure 1(e)), a reduction of one order of magnitude which is also confirmed by our experiments on real datasets, as reported in Figure 2(a). Moreover, in order to detect all inference channels holding in $\mathcal{F}(\mathcal{D}, \sigma)$, we can limit ourselves to retrieve only the inference channels in $MCh(k, Cl(\mathcal{D}, \sigma))$, thus taking in input $Cl(\mathcal{D}, \sigma)$ instead of $\mathcal{F}(\mathcal{D}, \sigma)$ and thus performing a much smaller number of checks. Algorithm 2 exploits the anti-monotonicity of frequency (Prop. 1) and the property of maximal inference channels (Prop. 2) to compute $MCh(k, Cl(\mathcal{D}, \sigma))$ from $Cl(\mathcal{D}, \sigma)$. Thanks to these two properties, Algorithm 2 is much faster, dramatically outperforming the naive inference channel detector (Algorithm 1), and scaling well even for very low support thresholds, as reported in Figure 2(b).

Algorithm 2 Optimized Inference Channel Detector

Input: $Cl(\mathcal{D}, \sigma), k$
Output: $MCh(k, Cl(\mathcal{D}, \sigma))$

- 1: $M = \{I \in Cl(\mathcal{D}, \sigma) \mid I \text{ is maximal}\}$;
- 2: $MCh(k, Cl(\mathcal{D}, \sigma)) = \emptyset$;
- 3: **for all** $J \in M$ **do**
- 4: **for all** $I \in Cl(\mathcal{D}, \sigma)$ **such that** $I \subseteq J$ **do**
- 5: **compute** f_I^J ;
- 6: **if** $0 < f_I^J < k$ **then**
- 7: **insert** $\langle C_I^J, f_I^J \rangle$ **in** $MCh(k, Cl(\mathcal{D}, \sigma))$;

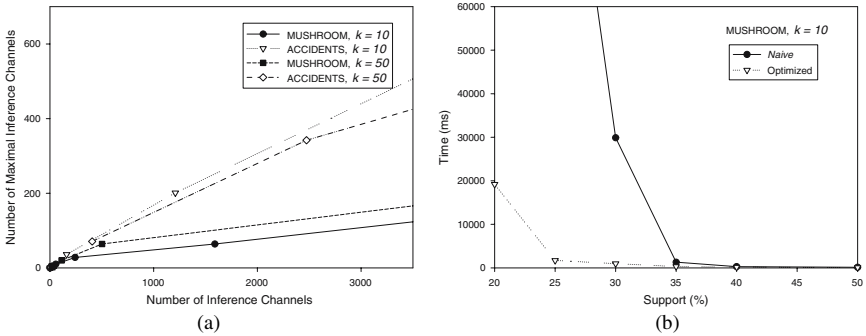


Fig. 2. Benefits of the condensed representation: size of the representations (a), and run time (b)

5 Anonymity vs. Accuracy: Empirical Observations

Algorithm 2 represents an optimized way to identify all threats to anonymity. Its performance revealed adequate in all our empirical evaluations using various datasets from the FIMI repository²; in all such cases the time improvement from the Naïve (Algorithm 1)

² <http://fimi.cs.helsinki.fi/data/>

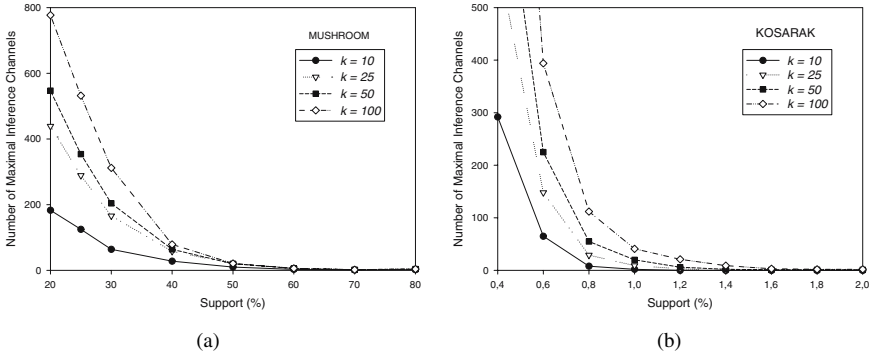


Fig. 3. Experimental results on cardinality of $\mathcal{MCh}(k, Cl(\mathcal{D}, \sigma))$ on two datasets

to the optimized algorithm is about one order of magnitude. This level of efficiency allows an interactive-iterative use of the algorithm by the analyst, aimed at finding the best trade-off among privacy and accuracy of the collection of patterns. To be more precise, there is a conflict among keeping the support threshold as low as possible, in order to mine all interesting patterns, and avoiding the generation of anonymity threats. The best solution to this problem is precisely to find out the minimum support threshold that generates a collection of patterns with no threats. The plots in Figure 3 illustrate this point: on the x -axis we report the minimum support threshold, on the y -axis we report the total number of threats (the cardinality of $\mathcal{MCh}(k, Cl(\mathcal{D}, \sigma))$), and the various curves indicate such number according to different values of the anonymity threshold k . In Figure 3(a) we report the plot for the MUSHROOM dataset (a dense one), while in Figure 3(b) we report the plot for the KOSARAK dataset which is sparse. In both cases, it is evident the value of the minimum support threshold that represents the best trade-off, for any given value of k . However, in certain cases, the best support threshold can still be too high to mine a sufficient quantity of interesting patterns. In such cases, the only option is to allow lower support thresholds and then to block the inference channels in the mining outcome. This problem, as stated before, is not covered in this paper for lack of space, and will be presented in a forthcoming paper.

6 Conclusions

We introduced in this paper the notion of k -anonymous patterns. Such notion serves as a basis for a formal account of the intuition that a collection of patterns, obtained by data mining techniques and made available to the public, should not offer any possibilities to violate the privacy of the individuals whose data are stored in the source database. To the above aim, we formalized the threats to anonymity by means of inference channel through frequent itemsets, and provided practical algorithms to detect such channels.

Other issues, emerging from our approach, are worth a deeper investigation and are left to future research. These include: (i) a thorough comparison of the various dif-

ferent approaches that may be used to block inference channels; (ii) a comprehensive empirical evaluation of our approach: to this purpose we are conducting a large-scale experiment with real life bio-medical data about patients to assess both applicability and scalability of the approach in a realistic, challenging domain; (iii) an investigation whether the proposed notion of privacy-preserving pattern discovery may be generalized to other forms of patterns and models.

In any case, the importance of the advocated form of privacy-preserving pattern discovery is evident: demonstrably trustworthy data mining techniques may open up tremendous opportunities for new knowledge-based applications of public utility and large societal and economic impact.

References

1. D. Agrawal and C. C. Aggarwal. On the design and quantification of privacy preserving data mining algorithms. In *Proceedings of the twentieth ACM PODS*, 2001.
2. R. Agrawal, T. Imielinski, and A. N. Swami. Mining association rules between sets of items in large databases. In *Proceedings of the 1993 ACM SIGMOD*.
3. R. Agrawal and R. Srikant. Privacy-preserving data mining. In *Proceedings of the 2000 ACM SIGMOD on Management of Data*.
4. R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules in Large Databases. In *Proceedings of the Twentieth VLDB*, 1994.
5. M. Atzori, F. Bonchi, F. Giannotti, and D. Pedreschi. k -anonymous patterns. Technical Report 2005-TR-17, ISTI - C.N.R., 2005.
6. T. Calders and B. Goethals. Mining all non-derivable frequent itemsets. In *Proceedings of the 6th PKDD*, 2002.
7. P. Fule and J. F. Roddick. Detecting privacy and ethical sensitivity in data mining results. In *Proc. of the 27th conference on Australasian computer science*, 2004.
8. D. Hand, H. Mannila, and P. Smyh. *Principles of Data Mining*. The MIT Press, 2001.
9. M. Kantarcioglu, J. Jin, and C. Clifton. When do data mining results violate privacy? In *Proceedings of the tenth ACM SIGKDD*, 2004.
10. D. Knuth. *Fundamental Algorithms*. Addison-Wesley, Reading, Massachusetts, 1997.
11. S. R. M. Oliveira, O. R. Zaiane, and Y. Saygin. Secure association rule sharing. In *Proc. of the 8th PAKDD*, 2004.
12. N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. In *Proc. ICDT '99*, 1999.
13. J. Pei, J. Han, and J. Wang. Closet+: Searching for the best strategies for mining frequent closed itemsets. In *SIGKDD '03*, 2003.
14. L. Sweeney. k -anonymity: a model for protecting privacy. *International Journal on Uncertainty Fuzziness and Knowledge-based Systems*, 10(5), 2002.
15. L. Sweeney. k -anonymity privacy protection using generalization and suppression. *International Journal on Uncertainty Fuzziness and Knowledge-based Systems*, 10(5), 2002.
16. V. S. Verykios, E. Bertino, I. N. Fovino, L. P. Provenza, Y. Saygin, and Y. Theodoridis. State-of-the-art in privacy preserving data mining. *SIGMOD Rec.*, 33(1):50–57, 2004.
17. M. J. Zaki and C.-J. Hsiao. Charm: An efficient algorithm for closed itemsets mining. In *2nd SIAM International Conference on Data Mining*, 2002.

Interestingness is Not a Dichotomy: Introducing Softness in Constrained Pattern Mining

Stefano Bistarelli^{1,2} and Francesco Bonchi³

¹ Dipartimento di Scienze, Università degli Studi “G. D’Annunzio”, Pescara, Italy

² Istituto di Informatica e Telematica, CNR, Pisa, Italy

³ Pisa KDD Laboratory, ISTI - C.N.R., Pisa, Italy

bista@sci.unich.it, francesco.bonchi@isti.cnr.it

Abstract. The paradigm of pattern discovery based on constraints was introduced with the aim of providing to the user a tool to drive the discovery process towards potentially *interesting* patterns, with the positive side effect of achieving a more efficient computation. So far the research on this paradigm has mainly focussed on the latter aspect: the development of efficient algorithms for the evaluation of constraint-based mining queries. Due to the lack of research on methodological issues, the constraint-based pattern mining framework still suffers from many problems which limit its practical relevance. As a solution, in this paper we introduce the new paradigm of pattern discovery based on *Soft Constraints*. Albeit simple, the proposed paradigm overcomes all the major methodological drawbacks of the classical constraint-based paradigm, representing an important step further towards practical pattern discovery.

1 Background and Motivations

During the last decade a lot of researchers have focussed their (mainly algorithmic) investigations on the computational problem of *Frequent Pattern Discovery*, i.e. mining patterns which satisfy a user-defined constraint of minimum frequency [1].

The simplest form of a frequent pattern is the frequent itemset.

Definition 1 (Frequent Itemset Mining). Let $\mathcal{I} = \{x_1, \dots, x_n\}$ be a set of distinct items, where an item is an object with some predefined attributes (e.g., price, type, etc.). An itemset X is a non-empty subset of \mathcal{I} . A transaction database \mathcal{D} is a bag of itemsets $t \in 2^{\mathcal{I}}$, usually called transactions. The support of an itemset X in database \mathcal{D} , denoted $\text{supp}_{\mathcal{D}}(X)$, is the number of transactions which are superset of X . Given a user-defined minimum support σ , an itemset X is called frequent in \mathcal{D} if $\text{supp}_{\mathcal{D}}(X) \geq \sigma$. This defines the minimum frequency constraint: $\mathcal{C}_{\text{freq}[\mathcal{D}, \sigma]}(X) \Leftrightarrow \text{supp}_{\mathcal{D}}(X) \geq \sigma$.

Recently the research community has turned its attention to more complex kinds of frequent patterns extracted from more structured data: *sequences*, *trees*, and *graphs*. All these different kinds of pattern have different peculiarities and application fields, but they all share the same computational aspects: a usually very large input, an exponential search space, and a too large solution set. This situation – too many data yielding too many patterns – is harmful for two reasons. First, performance degrades: mining generally becomes inefficient or, often, simply unfeasible. Second, the identification of the

fragments of interesting knowledge, blurred within a huge quantity of mostly useless patterns, is difficult. The paradigm of *constraint-based pattern mining* was introduced as a solution to both these problems. In such paradigm, it is the user which specifies to the system what is *interesting* for the current application: constraints are a tool to drive the mining process towards potentially interesting patterns, moreover they can be pushed deep inside the mining algorithm in order to fight the exponential search space curse, and to achieve better performance [15,20,25].

When instantiated to the pattern class of itemsets, the constraint-based pattern mining problem is defined as follows.

Definition 2 (Constrained Frequent Itemset Mining). *A constraint on itemsets is a function $\mathcal{C} : 2^X \rightarrow \{\text{true}, \text{false}\}$. We say that an itemset I satisfies a constraint if and only if $\mathcal{C}(I) = \text{true}$. We define the theory of a constraint as the set of itemsets which satisfy the constraint: $\text{th}(\mathcal{C}) = \{X \in 2^X \mid \mathcal{C}(X)\}$. Thus with this notation, the frequent itemsets mining problem requires to compute the set of all frequent itemsets $\text{th}(\mathcal{C}_{\text{freq}[\mathcal{D}, \sigma]})$. In general, given a conjunction of constraints \mathcal{C} the constrained frequent itemsets mining problem requires to compute $\text{th}(\mathcal{C}_{\text{freq}}) \cap \text{th}(\mathcal{C})$.*

Example 1. The following is an example mining query:

$$\mathcal{Q} : \text{th}(\mathcal{C}_{\text{freq}[\mathcal{D}, \sigma]}) \cap \text{th}(\mathcal{C}) \text{ where } \mathcal{C} = \{ \text{supp}(X) \geq 1500 \wedge \text{avg}(X.\text{weight}) \leq 5 \wedge \text{sum}(X.\text{price}) \geq 20 \}$$

It requires to mine, from database \mathcal{D} , all patterns which are frequent (have a support larger than 1500), have average weight less than 5 and a sum of prices greater than 20.

So far constraint-based frequent pattern mining has been seen as a query optimization problem, i.e., developing efficient, sound and complete evaluation strategies for constraint-based mining queries. Or in other terms, designing efficient algorithms to mine all and only the patterns in $\text{th}(\mathcal{C}_{\text{freq}}) \cap \text{th}(\mathcal{C})$. To this aim, properties of constraints have been studied comprehensively, and on the basis of such properties (e.g., anti-monotonicity, succinctness [20,18], monotonicity [11,17,6], convertibility [22], loose anti-monotonicity [9]), efficient computational strategies have been defined. Despite such effort, the constraint-based pattern mining framework still suffers from many problems which limit its practical relevance.

First of all, consider the example mining query \mathcal{Q} given above: *where do the three thresholds (i.e., 1500, 5 and 20) come from?* In some cases they can be precisely imposed by the application, but this is rarely the case. In most of the cases, they come from an exploratory mining process, where they are iteratively adjusted until a solution set of reasonable size is produced. This practical way of proceeding is in contrast with the basic philosophy of the constraint-based paradigm: constraints should represent what is a priori interesting, given the application background knowledge, rather than be adjusted accordingly to a preconceived output size. Another major drawback of the constraint-based pattern mining paradigm is its rigidity. Consider, for instance, the following three patterns (we use the notation $\langle v_1, v_2, v_3 \rangle$ to denote the three values corresponding to the three constraints in the conjunction in the example query \mathcal{Q}): $p_1 : \langle 1700, 0.8, 19 \rangle$, $p_2 : \langle 1550, 4.8, 54 \rangle$, and $p_3 : \langle 1550, 2.2, 26 \rangle$. The first pattern, p_1 , largely satisfies two out of the three given constraints, while slightly violates the third one. According to

the classical constraint-based pattern mining paradigm p_1 would be discarded as non interesting. Is such a pattern really *less interesting* than p_2 and p_3 which satisfy all the three constraints, but which are much less frequent than p_1 ? Moreover, is it reasonable, in real-world applications, that all constraints are equally important?

All these problems flow out from the same source: the fact that in the classical constraint-based mining framework, a constraint is a function which returns a boolean value $C : 2^X \rightarrow \{0, 1\}$. Indeed, *interestingness is not a dichotomy*.

This consideration suggests us a simple solution to overcome all the main drawbacks of constraint-based paradigm.

Paper Contributions and Organization

In this paper, as a mean to handle interestingness [26,16,24], we introduce the *soft constraint based pattern mining* paradigm, where constraints are no longer rigid boolean functions, but are “soft” functions, i.e., functions with value in a set A , which represents the set of interest levels or costs assigned to each pattern.

- The proposed paradigm is not rigid: a potentially interesting pattern is not discarded for just a slight violation of a constraint.
- Our paradigm creates an order of patterns w.r.t. interestingness (level of constraints satisfaction): this allows to say that a pattern is *more interesting* than another, instead of strictly dividing patterns in interesting and not interesting.
- From the previous point it follows that our paradigm allows to express *top-k* queries based on constraints: the data analyst can ask for the top-10 patterns w.r.t. a given description (a conjunction of soft constraints).
- Alternatively, we can ask to the system to return all and only the patterns which exhibit an interest level larger than a given threshold λ .
- The proposed paradigm allows to assign different weights to different constraints, while in the classical constraint-based pattern discovery paradigm all constraints were equally important.
- Last but not least, our idea is very simple and thus very general: it can be instantiated to different classes of patterns such as itemsets, sequences, trees or graphs.

For the reasons listed above, we believe that the proposed paradigm represents an important step further towards practical pattern discovery.

A nice feature of our proposal is that, by adopting the soft constraint based paradigm, we do not reject all research results obtained in the classical constraint-based paradigm; on the contrary, we fully exploit such algorithmic results. In other terms, our proposal is merely methodological, and it exploits previous research results that were mainly computational.

The paper is organized as follows. In the next Section we briefly review the theory of soft constraints and we define the soft constraint based pattern mining paradigm. In Section 3 we discuss possible alternative instances of the paradigm. In Section 4 we formally define the Soft Constraint Based Pattern Discovery paradigm. We then focus on one of the many possible instances of the proposed paradigm, and we implement it in a concrete Pattern Discovery System. Such a system is built as a wrapper around a classical constraint pattern mining system.

2 Introducing Soft Constraints

Constraint Solving is an emerging software technology for declarative description and effective solving of large problems. Many real life systems, ranging from network management [14] to complex scheduling [2], are analyzed and solved using constraint related technologies. The constraint programming process consists of the generation of requirements (constraints) and solution of these requirements, by specialized constraint solvers. When the requirements of a problem are expressed as a collection of boolean predicates over variables, we obtain what is called the *crisp* (or classical) Constraint Satisfaction Problem (CSP). In this case the problem is solved by finding any assignment of the variables that satisfies all the constraints.

Sometimes, when a deeper analysis of a problem is required, *soft constraints* are used instead. Several formalizations of the concept of soft constraints are currently available. In the following, we refer to the formalization based on *c-semirings* [5]: a semiring-based constraint assigns to each instantiation of its variables an associated value from a partially ordered set. When dealing with crisp constraints, the values are the boolean $\{0, 1\}$ representing the admissible and/or non-admissible values; when dealing with soft constraints the values are interpreted as preferences/costs. The framework must also handle the combination of constraints. To do this one must take into account such additional values, and thus the formalism must provide suitable operations for combination (\times) and comparison ($+$) of tuples of values and constraints. This is why this formalization is based on the mathematical concept of semiring.

Definition 3 (c-semirings [5,3]). A semiring is a tuple $\langle A, +, \times, \mathbf{0}, \mathbf{1} \rangle$ such that: A is a set and $\mathbf{0}, \mathbf{1} \in A$; $+$ is commutative, associative and $\mathbf{0}$ is its unit element; \times is associative, distributes over $+$, $\mathbf{1}$ is its unit element and $\mathbf{0}$ is its absorbing element. A *c-semiring* (“*c*” stands for “constraint-based”) is a semiring $\langle A, +, \times, \mathbf{0}, \mathbf{1} \rangle$ such that $+$ is idempotent with $\mathbf{1}$ as its absorbing element and \times is commutative.

Definition 4 (soft constraints [5,3]). Given a *c-semiring* $S = \langle A, +, \times, \mathbf{0}, \mathbf{1} \rangle$ and an ordered set of variables V over a finite domain D , a constraint is a function which, given an assignment $\eta : V \rightarrow D$ of the variables, returns a value of the *c-semiring*. By using this notation we define $C = \eta \rightarrow A$ as the set of all possible constraints that can be built starting from S , D and V .

In the following we will always use the word semiring as standing for *c-semiring*, and we will explain this very general concept by the point of view of pattern discovery.

Example 2. Consider again the mining query Q . In this context we have that the ordered set of variables V is $\langle \dots, D(X), avg(X.weight), sum(X.price) \rangle$, while the domain D is: $D(\dots, D(X)) = \mathbb{N}$, $D(avg(X.weight)) = \mathbb{R}^+$, and $D(sum(X.price)) = \mathbb{N}$. If we consider the classical *crisp* framework (i.e., hard constraints) we have the semiring $S_{Bool} = \langle \{true, false\}, \vee, \wedge, \dots, \dots \rangle$. A constraint C is a function $V \rightarrow D \rightarrow A$; for instance, $\dots, D(X) \rightarrow 1700 \rightarrow true$.

The $+$ operator is what we use to compare tuples of values (or patterns, in our context). Let us consider the relation \leq_S (where S stands for the specified semiring) over

A such that $a \leq_S b$ iff $a + b = b$. It is possible to prove that: \leq_S is a partial order; $+$ and \times are monotone on \leq_S ; $\mathbf{0}$ is its minimum and $\mathbf{1}$ its maximum, and $\langle A, \leq_S \rangle$ is a complete lattice with least upper bound operator $+$. In the context of pattern discovery $a \leq_S b$ means that the pattern b is *more interesting* than a , where interestingness is defined by a combination of soft constraints. When using (soft) constraints it is necessary to specify, via suitable combination operators, how the level of interest of a combination of constraints is obtained from the interest level of each constraint. The combined weight (or interest) of a combination of constraints is computed by using the operator $\otimes : \mathcal{C} \times \mathcal{C} \rightarrow \mathcal{C}$ defined as $(C_1 \otimes C_2)\eta = C_1\eta \times_S C_2\eta$.

Example 3. If we adopt the classical crisp framework, in the mining query \mathcal{Q} of Example 1 we have to combine the three constraints using the \wedge operator (which is the \times in the boolean semiring S_{Bool}). Consider for instance the pattern $p_1 : \langle 1700, 0.8, 19 \rangle$ for the ordered set of variables $V = \langle \cdot, \cdot, \cdot, \mathcal{D}(X), avg(X.weight), sum(X.price) \rangle$. The first and the second constraint are satisfied leading to the semiring level \cdot, \cdot , while the third one is not satisfied and has associated level \cdot, \cdot . Combining the three values with \wedge we obtain $\cdot, \cdot \wedge \cdot, \cdot \wedge \cdot, \cdot = \cdot, \cdot$ and we can conclude that the pattern $\langle 1700, 0.8, 19 \rangle$ is not interesting w.r.t. our purposes. Similarly, we can instead compute level \cdot, \cdot for pattern $p_3 : \langle 1550, 2.2, 26 \rangle$ corresponding to an interest w.r.t. our goals. Notice that using crisp constraints, the order between values only says that we are interested to patterns with semiring level \cdot, \cdot and not interested to patterns with semiring level \cdot, \cdot (that is semiring level $\cdot, \cdot \leq_{S_{Bool}} \cdot, \cdot$).

3 Instances of the Semiring

Dividing patterns in *interesting* and *non-interesting* is sometimes not meaningful nor useful. Most of the times we can say that each pattern is interesting with a specific level of preference. Soft constraints can deal with preferences by moving from the two values semiring S_{Bool} to other semirings able to give a finer distinction among patters (see [3] for a comprehensive guide to the semiring framework). For our scope the fuzzy and the weighted semiring are the most suitable.

Example 4 (fuzzy semiring). When using fuzzy semiring [12,23], to each pair constraint-pattern is assigned an interest level between 0 and 1, where 1 represents the best value (maximum interest) and 0 the worst one (minimum interest). Therefore the $+$ in this semiring is given by the *max* operator, and the order \leq_S is given by the usual \leq on real numbers. The value associated to a pattern is obtained by combining the constraints using the minimum operator among the semiring values. Therefore the \times in this semiring is given by the *min* operator. Recapitulating, the fuzzy semiring is given by $S_F = \langle [0, 1], max, min, 0, 1 \rangle$. The reason for such a max-min framework relies on the attempt to maximize the value of the least preferred tuple. Fuzzy soft constraints are able to model partial constraint satisfaction [13], so to get a solution even when the problem is overconstrained, and also prioritized constraints, that is, constraints with different levels of importance [10]. Figure 1 reports graphical representations of possible fuzzy instances of the constraints in \mathcal{Q} . Consider, for instance, the graphical representation of the frequency constraint in Figure 1(C_1). The dotted line describes the behavior

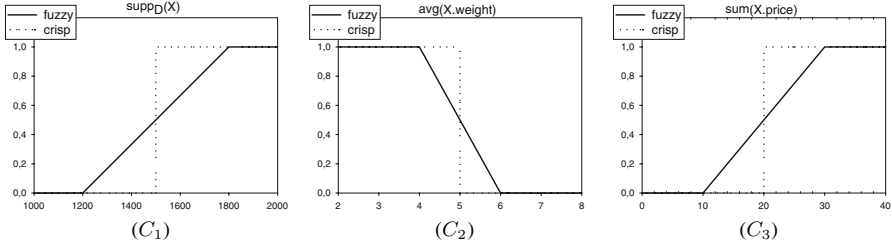


Fig. 1. Graphical representation of possible fuzzy instances of the constraints in \mathcal{Q}

of the *crisp* version (where 1 = *true* and 0 = *false*) of the frequency constraint, while the solid line describes a possible fuzzy instance of the same constraint. In this instance domain values smaller than 1200 yield 0 (uninteresting patterns); from 1200 to 1800 the interest level grows linearly reaching the maximum value of 1. Similarly for the other two constraints in Figure 1(C₂) and (C₃). In this situation for the pattern $p_1 = \langle 1700, 0.8, 19 \rangle$ we obtain that: $C_1(p_1) = 1$, $C_2(p_1) = 1$ and $C_3(p_1) = 0.45$. Since in the fuzzy semiring the combination operator \times is *min*, we got that the interest level of p_1 is 0.45. Similarly for p_2 and p_3 :

- $p_1 : C_1 \otimes C_2 \otimes C_3(1700, 0.8, 19) = \min(1, 1, 0.45) = 0.45$
- $p_2 : C_1 \otimes C_2 \otimes C_3(1550, 4.8, 54) = \min(1, 0.6, 1) = 0.6$
- $p_3 : C_1 \otimes C_2 \otimes C_3(1550, 2.2, 26) = \min(1, 1, 0.8) = 0.8$

Therefore, with this particular instance we got that $p_1 \leq_{S_F} p_2 \leq_{S_F} p_3$, i.e., p_3 is the most interesting pattern among the three.

Example 5 (weighted semiring). While fuzzy semiring associate a level of preference with each tuple in each constraint, in the weighted semiring tuples come with an associated cost. This allows one to model optimization problems where the goal is to minimize the total cost (time, space, number of resources, ...) of the proposed solution. Therefore, in the weighted semiring the cost function is defined by summing up the costs of all constraints. According to the informal description given above, the weighted semiring is $S_W = \langle \mathbb{R}^+, \min, \text{sum}, +\infty, 0 \rangle$. Consider, for instance, the graphical representation of the constraints in the query \mathcal{Q} in Figure 2. In this situation we got that:

- $p_1 : C_1 \otimes C_2 \otimes C_3(1700, 0.8, 19) = \text{sum}(50, 20, 205) = 275$
- $p_2 : C_1 \otimes C_2 \otimes C_3(1550, 4.8, 54) = \text{sum}(200, 120, 30) = 350$
- $p_3 : C_1 \otimes C_2 \otimes C_3(1550, 2.2, 26) = \text{sum}(200, 55, 190) = 445$

Therefore, with this particular instance we got that $p_3 \leq_{S_W} p_2 \leq_{S_W} p_1$ (remember that the order \leq_{S_W} correspond to the \geq on real numbers). In other terms, p_1 is the most interesting pattern w.r.t. this constraints instance.

The weighted and the fuzzy paradigm, can be seen as two different approaches to give a meaning to the notion of optimization. The two models correspond in fact to two definitions of social welfare in utility theory [19]: “*egalitarianism*”, which maximizes the minimal individual utility, and “*utilitarianism*”, which maximizes the sum of the

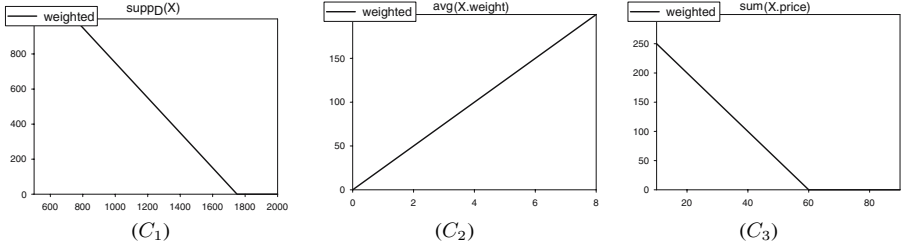


Fig. 2. Graphical representation of possible weighted instances of the constraints in \mathcal{Q}

individual utilities. The fuzzy paradigm has an egalitarianistic approach, aimed at maximizing the overall level of interest while balancing the levels of all constraints; while the weighted paradigm has an utilitarianistic approach, aimed at getting the minimum cost globally, even though some constraints may be neglected presenting a big cost. We believe that both approaches present advantages and drawbacks, and may be preferred to the other one depending on the application domain. Beyond the fuzzy and the weighted, many other possible instances of the semiring exist, and could be useful in particular applications. Moreover, it is worth noting that the cartesian product of semirings is a semiring [5] and thus it is possible to use the framework also to deal with multicriteria pattern selection.

Finally, note that the soft constraint framework is very general, and could be instantiated not only to unary constraints (as we do in this paper) but also to binary and k -ary constraints (dealing with two or more variables). This could be useful to extend the soft constraint based paradigm to association rules with “2-var” constraints [18].

4 Soft Constraint Based Pattern Mining

In this Section we instantiate soft constraint theory to the pattern discovery framework.

Definition 5 (Soft Constraint Based Pattern Mining). *Let \mathcal{P} denote the domain of possible patterns. A soft constraint on patterns is a function $\mathcal{C} : \mathcal{P} \rightarrow A$ where A is the carrier set of a semiring $S = \langle A, +, \times, \mathbf{0}, \mathbf{1} \rangle$. Given a combination of soft constraints $\otimes \mathcal{C}$, we define two different problems:*

λ -interesting: *given a minimum interest threshold $\lambda \in A$, it is required to mine the set of all λ -interesting patterns, i.e., $\{p \in \mathcal{P} \mid \otimes \mathcal{C}(p) \geq \lambda\}$.*

top- k : *given a threshold $k \in \mathbb{N}$, it is required to mine the top- k patterns $p \in \mathcal{P}$ w.r.t. the order \leq_S .*

Note that the Soft Constraint Based Pattern Mining paradigm just defined, has many degrees of freedom. In particular, it can be instantiated: (i) on the domain of patterns \mathcal{P} in analysis (e.g., itemsets, sequences, trees or graphs), (ii) on the semiring $S = \langle A, +, \times, \mathbf{0}, \mathbf{1} \rangle$ (e.g., fuzzy, weighted or probabilistic), and (iii) on one of the two possible mining problems, i.e., λ -interesting or top- k mining.

In the rest of this paper we will focus on concretizing a simple instance of this very general paradigm: λ -interesting_{fuzzy} on the pattern class of itemsets.

4.1 Mining λ -Interesting Itemsets on the Fuzzy Semiring

Definition 6. Let $\mathcal{I} = \{x_1, \dots, x_n\}$ be a set of items, where an item is an object with some predefined attributes (e.g., price, type, etc.). A soft constraint on itemsets, based on the fuzzy semiring, is a function $\mathcal{C} : 2^{\mathcal{I}} \rightarrow [0, 1]$. Given a combination of such soft constraints $\otimes \mathcal{C} \equiv \mathcal{C}_1 \otimes \dots \otimes \mathcal{C}_n$, we define the interest level of an itemset $X \in 2^{\mathcal{I}}$ as $\otimes \mathcal{C}(X) = \min(\mathcal{C}_1(X), \dots, \mathcal{C}_n(X))$. Given a minimum interest threshold $\lambda \in]0, 1]$, the λ -interesting itemsets mining problem, requires to compute $\{X \in 2^{\mathcal{I}} \mid \otimes \mathcal{C}(X) \geq \lambda\}$.

In the following we describe how to build a concrete *pattern discovery system* for λ -interesting_{fuzzy} itemsets mining, as a wrapper around a classical constraint pattern mining system. The basic components which we use to build our system are the following:

A crisp constraints solver - i.e., a system for mining constrained frequent itemsets, where constraints are classical binary functions, and not soft constraints. Or in other terms, a system for solving the problem in Definition 2. To this purpose we adopt the system which we have developed at Pisa KDD Laboratory within the *P³D* project¹. Such a system is a general Apriori-like algorithm which, by means of *data reduction* and *search space pruning*, is able to push a wide variety of constraints (practically all possible kinds of constraints which have been studied and characterized so far [9]) into the frequent itemsets computation. Based on the algorithmic results developed in the last years by our lab (e.g., [6,7,8,9,21]), our system is very efficient and robust, and to our knowledge, is the unique existing implementation of this kind.

A language of constraints - to express, by means of queries containing conjunctions of constraints, what is interesting for the given application. The wide repertoire of constraints that we admit, comprehends the frequency constraint $(\cdot, \cdot, \mathcal{D}(X) \geq \sigma)$, and all constraints defined over the following aggregates²: *min, max, count, sum, range, avg, var, median, std, md*.

A methodology to define the interest level - that must be assigned to each pair itemset-constraint. In other terms, we need to provide the analyst with a simple methodology to define how to assign for each constraint and each itemset a value in the interval $[0, 1]$, as done, for instance, by the graphical representations of constraints in Figure 1. This methodology should provide the analyst with a knob to adjust the *softness level* of each constraint in the conjunction, and a knob to set the *importance* of each constraint in the conjunction.

Let us focus on the last point. Essentially we must describe how the user can define the fuzzy behavior of a soft constraint. We restrict our system to constraints which behave as those ones in Figure 1: they return a value which grows linearly from 0 to 1 in a certain interval, while they are null before the interval and equal to 1 after the interval. To describe such a simple behavior we just need two parameters: a value associated to the center of the interval (corresponding to the 0.5 fuzzy semiring value), and a parameter to adjust the width of the interval (and consequently the gradient of the function).

¹ <http://www-kdd.isti.cnr.it/p3d/index.html>

² *Range* is $(\max - \min)$, *var* is for variance, *std* is for standard deviation, *md* is for mean deviation.

Definition 7. A soft constraint \mathcal{C} on itemsets, based on the fuzzy semiring, is defined by a quintuple $\langle \text{Agg}, \text{Att}, \theta, t, \alpha \rangle$, where:

- $\text{Agg} \in \{\text{supp}, \text{min}, \text{max}, \text{count}, \text{sum}, \text{range}, \text{avg}, \text{var}, \text{median}, \text{std}, \text{md}\}$;
- Att is the name of the attribute on which the aggregate agg is computed (or the transaction database, in the case of the frequency constraint);
- $\theta \in \{\leq, \geq\}$;
- $t \in \mathbb{R}$ corresponds to the center of the interval and it is associated to the semiring value 0.5;
- $\alpha \in \mathbb{R}^+$ is the softness parameter, which defines the inclination of the preference function (and thus the width of the interval).

In particular, if $\theta = \leq$ (as in Figure 1(C_2)) then $\mathcal{C}(X)$ is 1 for $X \leq (t - \alpha t)$, is 0 for $X \geq (t + \alpha t)$, and is linearly decreasing from 1 to 0 within the interval $[t - \alpha t, t + \alpha t]$. The other way around if $\theta = \geq$ (as, for instance, in Figure 1(C_3)). Note that if the softness parameter α is 0, then we obtain the crisp (or hard) version of the constraint.

Example 6. Consider again the query \mathcal{Q} given in Example 1, and its fuzzy instance graphically described by Figure 1. Such query can be expressed in our constraint language as: $\langle ., ., ., \mathcal{D}, \geq, 1500, 0.2 \rangle, \langle \text{avg}, \text{weight}, \leq, 5, 0.2 \rangle, \langle \text{sum}, \text{price}, \geq, 20, 0.5 \rangle$.

Since the combination operator \times in min , increasing the importance of a constraint w.r.t. the others in the combination means to force the constraint to return lower values for not really satisfactory patterns. By decreasing the softness parameter α , we increase the gradient of the function making the shape of the soft constraint closer to a crisp constraint. This translates in a better value for patterns X which were already behaving well w.r.t. such constraint ($\mathcal{C}(X) > 0.5$), and in a lower value for patterns which were behaving not so well ($\mathcal{C}(X) < 0.5$). Decreasing the gradient (increasing α) instead means to lower the importance of the constraint itself: satisfying or not satisfying the constraint does not result in a big fuzzy value difference. Additionally, by operating on t , we can increase the “severity” of the constraint w.r.t. those patterns which were behaving not so well. Therefore, the knob to increase or decrease the importance of a constraint is not explicitly given, because its role, in the fuzzy semiring, can be played by a combined action on the two knobs α and t .

Example 7. Consider again the query \mathcal{Q} given in Example 1, and its fuzzy instance: $\langle ., ., ., \mathcal{D}, \geq, 1500, 0.2 \rangle, \langle \text{avg}, \text{weight}, \leq, 5, 0.2 \rangle, \langle \text{sum}, \text{price}, \geq, 20, 0.5 \rangle$. As we stated in Example 4, it holds that $p_2 \leq_{SF} p_3$. In particular, p_2 is better than p_3 w.r.t. constraint C_3 , while p_3 is better than p_2 w.r.t. constraint C_2 . Suppose now that we increase the importance of C_3 , e.g., $\langle \text{sum}, \text{price}, \geq, 28, 0.25 \rangle$. We obtain that $p_3 \leq_{SF} p_2$:

- $p_2 : C_1 \otimes C_2 \otimes C_3(1550, 4.8, 54) = \min(1, 0.6, 1) = 0.6$
- $p_3 : C_1 \otimes C_2 \otimes C_3(1550, 2.2, 26) = \min(1, 1, 0.35) = 0.35$

In [5,4] it has been proved that, when dealing with the fuzzy framework, computing all the solution better than a threshold λ can be performed by solving a crisp problem where all the constraint instances with semiring level lower than λ have been assigned level λ , and all the instances with semiring level greater or equal to λ have been assigned level 1. Using this theoretical result, and some simple arithmetic we can transform each soft constraint in a corresponding crisp constraint.

Definition 8. Given a fuzzy soft constraint $\mathcal{C} \equiv \langle \text{Agg}, \text{Att}, \theta, t, \alpha \rangle$, and a minimum interest threshold λ , we define the crisp translation of \mathcal{C} w.r.t. λ as:

$$C_{crisp}^\lambda \equiv \begin{cases} \text{Agg}(\text{Att}) \geq t - \alpha t + 2\lambda\alpha t, & \text{if } \theta = \geq \\ \text{Agg}(\text{Att}) \leq t + \alpha t - 2\lambda\alpha t, & \text{if } \theta = \leq \end{cases}$$

Example 8. The crisp translation of the soft constraint $\langle \text{sum}, \text{price}, \geq, 20, 0.5 \rangle$ is $\text{sum}(X.\text{price}) \geq 26$ for $\lambda = 0.8$, while it is $\text{sum}(X.\text{price}) \geq 18$ for $\lambda = 0.4$.

Proposition 1. Given the vocabulary of items \mathcal{I} , a combination of soft constraints $\otimes \mathcal{C} \equiv \mathcal{C}_1 \otimes \dots \otimes \mathcal{C}_n$, and a minimum interest threshold λ . Let \mathcal{C}' be the conjunction of crisp constraints obtained by conjoining the crisp translation of each constraint in $\otimes \mathcal{C}$ w.r.t. λ : $\mathcal{C}' \equiv \mathcal{C}_1^\lambda_{crisp} \wedge \dots \wedge \mathcal{C}_n^\lambda_{crisp}$. It holds that: $\{X \in 2^{\mathcal{I}} \mid \otimes \mathcal{C}(X) \geq \lambda\} = \dots (\mathcal{C}')$.

Proof (sketch). The soundness of the mapping come from the result in [5]. We here have to only give a justification of the formula in Definition 8. This is done by means of Figure 3(b), that shows a graphical representation of the simple arithmetic problem and its solutions.

Therefore, if we adopt the fuzzy semiring, we can fully exploit a classical constraint-based pattern discovery system (and all algorithmic results behind it), by means of a simple translation from soft to crisp constraints. This is exactly what we have done, obtaining a pattern discovery system based on soft constraints built as a wrapper around a classical constraint-based mining system.

4.2 Experimental Analysis

We have conducted some experiments in order to asses the concrete effects obtained by manipulating the α , t and λ parameters. To this purpose we have compared 5 different instances (described in Figure 3(a)) of the query \mathcal{Q} :

$$\langle \dots, \mathcal{D}, \geq, t, \alpha \rangle \langle \text{avg}, \text{weight}, \leq, t, \alpha \rangle, \langle \text{sum}, \text{price}, \geq, t, \alpha \rangle$$

where the transactional dataset \mathcal{D} , is the well known RETAIL dataset, donated by Tom Brijs and contains the (anonymized) retail market basket data from an anonymous Belgian retail store³; and the two attributes *weight* and *price* have been randomly generated with a gaussian distribution within the range $[0, 150000]$.

Figure 3(c) reports the number of solutions for the given five queries at different λ thresholds. Obviously as λ increases the number of solutions shrinks accordingly. This behavior is also reflected in queries evaluation times, reported in Figure 3(d): the bigger is the size of the solution set, the longer is the associated computation.

Comparing queries \mathcal{Q}_1 , \mathcal{Q}_2 and \mathcal{Q}_3 , we can gain more insight about the α parameter. In fact, the three queries differ only by the α associated with one constraint (the frequency constraint). We can observe that, if the λ threshold is not too much selective, increasing the α parameter (i.e., the size of the soft interval), the number of solutions grows. Notice however that, when λ becomes selective enough (i.e., $\lambda > 0.5$), increasing the softness parameter we obtain an opposite behavior. This is due to the fact that, if on one hand a more soft constraint is less severe with patterns not good enough, on the other hand it is less generous with good patterns, which risk to be discarded by an high λ threshold.

³ <http://fimi.cs.helsinki.fi/data/>

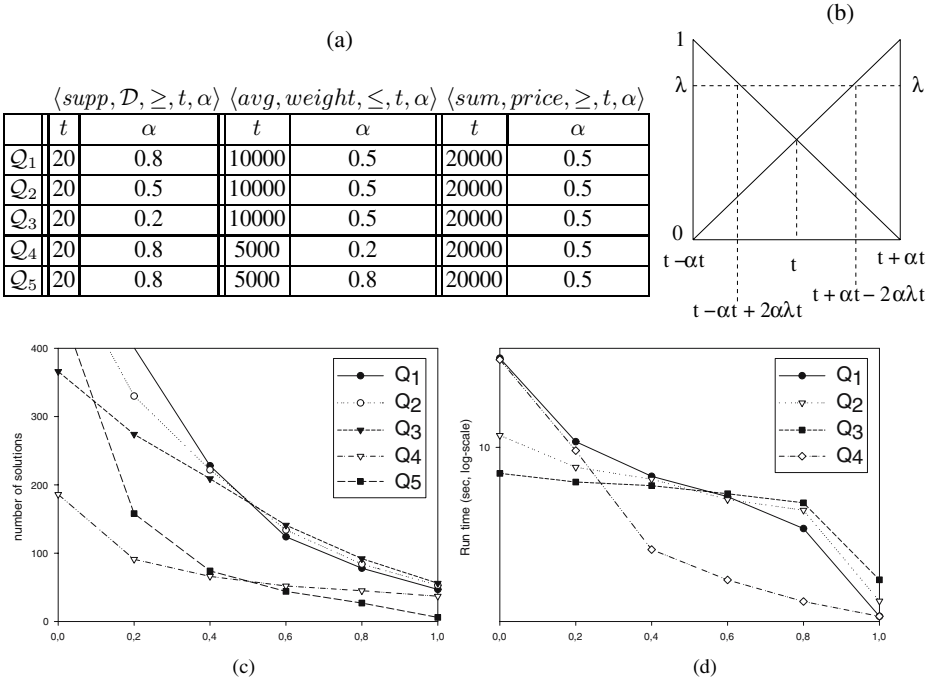


Fig. 3. (a) description of queries experimented, (b) graphical proof to Proposition 1, (c) and (d) experimental results with λ ranging in $[0, 1]$

References

1. R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules in Large Databases. In *Proceedings of the Twentieth International Conference on Very Large Databases (VLDB'94)*, 1994.
2. J. Bellone, A. Chamard, and C. Pradelles. Plane - an evolutive planning system for aircraft production. In *Proc. 1st International Conference on Practical Applications of Prolog (PAP92)*, 1992.
3. S. Bistarelli. *Semirings for Soft Constraint Solving and Programming*, volume 2962 of *Lecture Notes in Computer Science*. Springer, 2004.
4. S. Bistarelli, P. Codognet, and F. Rossi. Abstracting soft constraints: Framework, properties, examples. *Artificial Intelligence*, (139):175–211, July 2002.
5. S. Bistarelli, U. Montanari, and F. Rossi. Semiring-based Constraint Solving and Optimization. *Journal of the ACM*, 44(2):201–236, Mar 1997.
6. F. Bonchi, F. Giannotti, A. Mazzanti, and D. Pedreschi. ExAMiner: Optimized level-wise frequent pattern mining with monotone constraints. In *Proceedings of the Third IEEE International Conference on Data Mining (ICDM'03)*, 2003.
7. F. Bonchi, F. Giannotti, A. Mazzanti, and D. Pedreschi. ExAnte: Anticipated data reduction in constrained pattern mining. In *Proceedings of the 7th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD'03)*, 2003.
8. F. Bonchi and C. Lucchese. On closed constrained frequent pattern mining. In *Proceedings of the Fourth IEEE International Conference on Data Mining (ICDM'04)*, 2004.

9. F. Bonchi and C. Lucchese. Pushing tougher constraints in frequent pattern mining. In *Proceedings of the Ninth Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD'05)*, Hanoi, Vietnam, 2005.
10. A. Borning, M. Maher, A. Martindale, and M. Wilson. Constraint hierarchies and logic programming. In *Proc. 6th International Conference on Logic Programming*, 1989.
11. C. Bucila, J. Gehrke, D. Kifer, and W. White. DualMiner: A dual-pruning algorithm for itemsets with constraints. In *Proceedings of the 8th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD'02)*, 2002.
12. D. Dubois, H. Fargier, and H. Prade. The calculus of fuzzy restrictions as a basis for flexible constraint satisfaction. In *Proc. IEEE International Conference on Fuzzy Systems*, pages 1131–1136. IEEE, 1993.
13. E. Freuder and R. Wallace. Partial constraint satisfaction. *AI Journal*, 58, 1992.
14. T. Frühwirth and P. Brisset. Optimal planning of digital cordless telecommunication systems. In *Proc. PACT97*, London, UH, 1997.
15. J. Han, L. V. S. Lakshmanan, and R. T. Ng. Constraint-based, multidimensional data mining. *Computer*, 32(8):46–50, 1999.
16. R. Hilderman and H. Hamilton. *Knowledge Discovery and Measures of Interest*. Kluwer Academic, Boston, 2002.
17. S. Kramer, L. D. Raedt, and C. Helma. Molecular feature mining in hiv data. In *Proceedings of the 7th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD'01)*, 2001.
18. L. V. S. Lakshmanan, R. T. Ng, J. Han, and A. Pang. Optimization of constrained frequent set queries with 2-variable constraints. In *Proceedings of the ACM International Conference on Management of Data (SIGMOD'99)*, 1999.
19. H. Moulin. *Axioms for Cooperative Decision Making*. Cambridge University Press, 1988.
20. R. T. Ng, L. V. S. Lakshmanan, J. Han, and A. Pang. Exploratory mining and pruning optimizations of constrained associations rules. In *Proceedings of the ACM International Conference on Management of Data (SIGMOD'98)*, 1998.
21. S. Orlando, P. Palmerini, R. Perego, and F. Silvestri. Adaptive and Resource-Aware Mining of Frequent Sets. In *Proc. of the 2002 IEEE Int. Conference on Data Mining (ICDM'02)*, pages 338–345, Maebashi City, Japan, Dec. 2002.
22. J. Pei and J. Han. Can we push more constraints into frequent pattern mining? In *Proceedings of the 6th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD'00)*, 2000.
23. Z. Ruttkay. Fuzzy constraint satisfaction. In *Proc. 3rd IEEE International Conference on Fuzzy Systems*, pages 1263–1268, 1994.
24. S. Sahar. Interestingness via what is not interesting. In *Proc. of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD'99)*.
25. R. Srikant, Q. Vu, and R. Agrawal. Mining association rules with item constraints. In *Proceedings of the 3rd ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD'97)*, 1997.
26. P.-N. Tan, V. Kumar, and J. Srivastava. Selecting the right interestingness measure for association patterns. In *Proc. of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD'2002)*.

Generating Dynamic Higher-Order Markov Models in Web Usage Mining

João Borges¹ and Maël Levene²

¹ School of Engineering, University of Porto,
R. Dr. Roberto Frias, 4200 - Porto, Portugal
jlborges@fe.up.pt

² School of Computer Science and Information Systems, Birkbeck,
University of London, Malet Street, London WC1E 7HX, UK
mlevene@dcs.bbk.ac.uk

Abstract. Markov models have been widely used for modelling users' web navigation behaviour. In previous work we have presented a dynamic clustering-based Markov model that accurately represents second-order transition probabilities given by a collection of navigation sessions. Herein, we propose a generalisation of the method that takes into account higher-order conditional probabilities. The method makes use of the state cloning concept together with a clustering technique to separate the navigation paths that reveal differences in the conditional probabilities. We report on experiments conducted with three real world data sets. The results show that some pages require a long history to understand the users choice of link, while others require only a short history. We also show that the number of additional states induced by the method can be controlled through a probability threshold parameter.

1 Introduction

Modelling the web navigation data is a challenging task having to gain insight into the behaviour of the web and internet-based applications. Data characterising web navigation can be collected from the web or client-based log files, enabling the reconstruction of the navigation session [15]. A session is usually defined as a sequence of pages viewed by a user within a given time window. The basic historical method to extract a user's navigation data has been called *session reconstruction*, and this method have been applied in several contexts including personalisation, link prediction, e-commerce analysis, adaptive web personalisation and web page re-engineering [10].

Several authors have proposed the use of Markov models to represent a collection of the web navigation session. Piotrowski et al. [12] proposed a method to induce the collection of long sequences based on the *longest common subsequence*, while Dehondt et al. [7] proposed a technique to build *kth-order* Markov models and combine the models to include the higher order models covering each page. On the other hand, Sarrail [13] presented a methodology showing how Markov models have been used in link prediction applications, while Zheng et al. [16] inferred a Markov model from the web navigation data to evaluate page co-citation and clustering quality.

An alternative method of modeling navigation session behavior is proposed by Schech et al. [14] using a hidden Markov model (HMM) to model the collection of a hidden Markov model (HMM) to model the next page accessed, while Donghan and Jny [8] proposed a hybrid-order Markov model to model web page access. In addition, Chen and Zhang [6] use a Prediction by Partial Matching (PPM) model to model the sequence of nodes.

In previous work, we proposed a model of web navigation session as a Hidden Markov Process (HMP) [1,2]. A HMP consists of a hidden Markov model, which is a sequence of the N -gram concept [5] to achieve increased accuracy by increasing the order of the Markov chain; for the full definition of the HMP concept see [2]. In [2] an algorithm is proposed for the efficient evaluation of the HMP model, and in [3] we have shown that the algorithm is computationally efficient in terms of space and time. In [4] we extended the HMP model with a dynamic clustering-based method for state cloning [9] to accurately model second-order conditional probabilities; the method is described in Section 2. In this work, we generalize the method given in [4] to higher-order conditional probabilities.

Modern web mining requires techniques such as clustering, association rule mining and sequential pattern mining to each form a part in navigation modeling [10], and do not take into account the order in which pages were accessed. This limitation has been addressed by building a sequence of higher-order Markov models with a method that chooses the best model order in each case [7]. However, we argue that a single model is sufficient for modeling the variable length history of page history.

The method we propose in Section 3 aims to allow the cloning of a state in a Markov model to allow modeling of page history and the choice of length of history. In this way the order of a given state reflects the n -order conditional probability of the next page. In addition, the proposed model maintains the fundamental properties of the HMP model [1], while providing a simple and efficient algorithm for modeling the navigation session history in an account of page view.

In Section 2 we review the formal definition of the dynamic clustering method, in Section 3 we extend the method to model higher-order probabilities, and in Section 4 we review the experimental results. Finally, in Section 5 we give our concluding remarks.

2 Background

In previous work [1,2] we proposed a model of web navigation data as a Hidden Markov Process (HMP), which consists of a hidden Markov model. We now review the HMP model with the aid of an example.

Consider a website with seven web pages, $\{A_1, A_2, \dots, A_7\}$, and the collection of navigation session given on the left side of Figure 1 (NOS). The sequence of occurrence of each session). A navigation session gives a sequence

ence of page viewed by a user within a given time window. To each webpage we associate a node in the model. In addition, the anchor page, S , represents the starting page of every navigation session, and the anchor page, F , represents the last page of every navigation session. The edges (transition) corresponding to each pair of page visited in the sequence, a transition from S to the starting page of a session, and a transition from the last page of a session to F . The model is incrementally built by processing the complete collection of navigation sessions.

Session	NOS
A_1, A_2, A_3	3
A_1, A_2, A_4	2
A_5, A_2, A_3	1
A_5, A_2, A_4	1
A_6, A_2, A_3	2
A_6, A_2, A_4	1
A_1, A_2, A_4, A_7	1
A_5, A_2, A_4, A_7	3
A_6, A_2, A_4, A_7	2

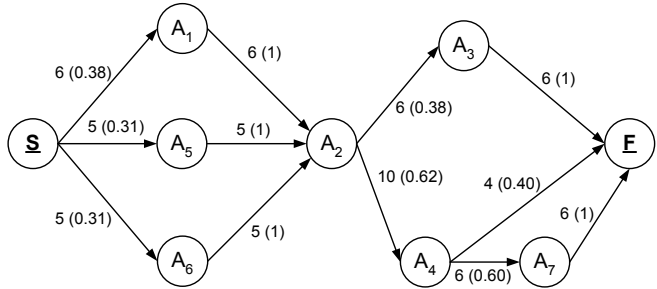


Fig. 1. A collection of user navigation sessions and the corresponding first-order model

A transition probability is defined by the ratio of the number of transitions was observed and the number of the anchor page was visited. The right side of Figure 1 shows a sequence of the first-order model corresponding to the navigation. Next, a linear model is given, the number of the linear was observed and the number in a given time give it is defined probability.

In [4] we proposed a method to increase the HPG precision in order to accurately represent second-order probabilities. The method is based on a cloning operation, where a page is duplicated if second-order probabilities diverge from the corresponding second-order probabilities. In addition, the method is a clustering algorithm to identify the best way to distribute a set of in-links between a page and its clones. We now present the essential properties of the model proposed in [4], which we extend to higher-order probabilities in Section 3.

Given a model with a set $\{S, A_1, \dots, A_n, F\}$, we let w_i represent the number of the page corresponding to A_i was visited, $w_{i,j}$ be the number of the in-link from A_i to A_j was observed, and $w_{i,j,k}$ be the number of the sequence A_i, A_j, A_k was observed. In addition, we let $p_{i,j} = w_{i,j}/w_i$ be the first-order transition probability from A_i to A_j , and $p_{i,k,j} = w_{i,k,j}/w_{i,k}$ be the second-order transition probability. Also, the accuracy holds, γ , is the highest admissible difference between a first-order and a second-order probability; a model is said to be γ -accurate if the linear model does not violate the constraint by γ .

In the example given in Figure 1, the sequence of a navigation behavior is linear has $p_{1,23} = p_{1,24} = 0.5$. Therefore, for $\gamma = 0.1$, the page A_2 is not accurate, since $|p_{1,23} - p_{2,3}| > 0.1$, and needs to be cloned. To clone a page A_2 , we let each in-link

denote a vector of second-order transition probabilities; each of the vectors' components corresponds to an ordered pair of nodes in A_2 . In the example, node A_2 has three in-links and two out-links, indicating the vector of second-order probabilities: for $i = \{3, 4\}$ we have $P_{1,2i} = \{0.5, 0.5\}$, $P_{5,2i} = \{0.2, 0.8\}$ and $P_{6,2i} = \{0.4, 0.6\}$.

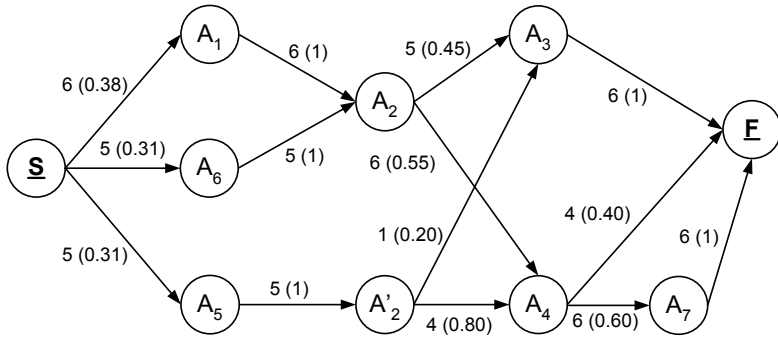


Fig. 2. The second-order HPG model obtained when applying the dynamic clustering method with $\gamma = 0.1$ to the first-order model given in Figure 1

The method also implements a k -mean clustering algorithm on the collection of second-order vectors, in order to identify groups of similar vectors with respect to γ . Figure 2 shows the result of applying the method to node A_2 . Since $|p_{1,2i} - p_{6,2i}| < 0.1$, for $i = \{3, 4\}$, links from A_1 and A_6 are assigned to one cluster, the link from A_5 is assigned to the other cluster. The transition counts for the out-links are updated as follows: $w_{2,3} = w_{1,2,3} + w_{6,2,3}$, $w_{2,4} = w_{1,2,4} + w_{6,2,4}$, $w'_{2,3} = w_{5,2,3}$, $w'_{2,4} = w_{5,2,4}$. Note that node A_4 is accessed, since all its in-links have an equivalent source node, and, moreover, every node having a single in-link is accessed by definition. Therefore, the model given in Figure 2 accurately represents every second-order transition probability.

3 A Dynamic Clustering Method to Model Higher-Order Probabilities

We now extend the method presented in [4] to incorporate higher-order probabilities. In a second-order HPG model, the transition probability for a given node is considered to be accessed, if all in-links of it indicate identical second-order probabilities. Similarly, in a third-order model every two-link path to a node is considered identical third-order probabilities. In general, to access a model n -order probabilities each $(n - 1)$ -length path to a node is considered identical n -order conditional probabilities. Each one of the n -order conditional probabilities are obtained from the n -gram counts.

In the following, we let the length of a path be denoted by the number of links it is composed of, and we call the length of the path for a node

o the age γ a e he \dots of hi \dots a e; $d = 0$ co, e ond o the age γ a e and $d = n - 1$ co, e ond o the fa, he \dots a e fo the age when a e ing he odel fo, o de, n . We le $w_{1,\dots,n}$, e e en he n -g, a co n, and $p_{i,\dots,j,kt} = w_{i,\dots,j,k,t}/w_{i,\dots,j,k}$, e e en he n -o, de, condi onal, obabili y of going o a e A_t given ha the $(n - 1)$ -leng h a h A_i, \dots, A_j, A_k wa followed. Al o, we le \vec{l} , e e en a a h and $p_{\vec{l},kt}$ he condi onal, obabili y of, an i ion (A_k, A_t) given he a h \vec{l} . Al o, we le $\vec{l}_{[d]}$ be he a e a de h d on \vec{l} and $v_{\vec{l}}$ be he vec o, of n -o, de, condi onal, obabili e given a h \vec{l} . If a e y need c_y clone, we le y_i , wi h $i = \{1, \dots, c_y\}$, e e en y and i $c_y - 1$ addi onal clone. Finally, we le \vec{l}_c be he cl, e, o which a h \vec{l} wa a, ign ed.

Fo, a a e x , he n -o, de, condi onal, obabili e a e a e ed in h ee, e, :

- (i) A ly a b, ea h, \dots , ea, ch, oced, e o ind ce he $(n - 1)$ -leng h in- a h o a e x , e i a e he co, e onding n -o, de, condi onal, obabili e and, fo, ea, ch, a h, \vec{l} , o e he condi onal, obabili e in a vec o, $v_{\vec{l}}$ (he vec o, ' di en ion i given by he n. be, of o -lin. fo x). If he diffe, ence be ween a condi onal, obabili y and he co, e onding, an i ion, obabili y i g ea e, han γ , label he a e a need ing o be cloned.
- (ii) If x need cloning, a ly he k - ean algo, i h o he, obabili y vec o, , $v_{\vec{l}}$. The n. be, of cl, e, k i inc, e en ed n il in he nal, ol ion, and in eve, y cl, e, , he di, stance be ween ea, ch, vec o, and i. cen, oid i. alle, han γ .
- (iii) Iden ify a e ha need o be cloned o e a a e he a h o x . Sa e incl ded in a h o x a e a e ed in de cend ing de h o, de, fo, $d = n - 1$ o $d = 0$. Fo, de h d , we le a a e x of a a h o x , who e la, a e i y , be na ed a y a h- e x o x . Th, o e a a e a h o x , a e y a de h, d , need a a any clone a he n. be, of di inc, a h, e x e wi h he a e leng h ha a e a, ign ed o diffe, en cl, e, . The weigh, of he in and o -lin. of y and i clone a e de e, ined by he n -g, a co n. Af e, cloning y he in- a h o x need o be da ed.

We now, e en an exa, ple of he e hod and a e do- code de c i ion. In a ic la, , we eval a e he hi, d- o, de, obabili e fo, he odel in Fig, e 2. The condi onal, obabili e ind ced by he a h o A_4 a e: fo, $i = \{7, F\}$ we have $p_{12,4i} = \{0.33, 0.67\}$, $p_{62,4i} = \{0.67, 0.33\}$ and $p_{52,4i} = \{0.75, 0.25\}$. Th, o e a e obabili e a e no clo e o he co, e onding, econd- o, de, obabili e, A_4 i no hi, d- o, de, acc, a e fo, $\gamma = 0.1$. Table 1 give he in- a h o A_4 , he hi, d- o, de, condi onal, obabili e and he e ling cl, e, ing a, ign en. A e l, , a e A_2 fo, $d = 1$ need one clone, and fo, $d = 0$, a e A_4 al o need one clone. Fig, e 3 give he e ling hi, d- o, de, odel.

In Fig, e 3, he a h S, A_1, A_2, A_4 ha obabili y e i a e of $0.38 \cdot 1.00 \cdot 0.50 = 0.19$. I can be e en ha in Fig, e 1, fo, a o al of 16 e ion, 3 begin wi h he 3- g, a A_1, A_2, A_4 e ling in a obabili y e i a e of 0.19. Al o, acco, ding o he hi, d- o, de, odel, a h S, A_5, A_2, A_4, A_7 ha obabili y $0.31 \cdot 1.00 \cdot 0.80 \cdot 0.71 = 0.18$. I can be e en ha in he in da a 3 e ion

Table 1. The paths to A_4 , the third-order conditional probabilities and the resulting clustering assignment

$d = 2$	$d = 1$	$d = 0$	3rd order vectors		cluster
A_1	A_2	A_4	0.33	0.67	1
A_6	A_2	A_4	0.67	0.33	2
A_5	A_2	A_4	0.75	0.25	2

begin with A_5, A_2, A_4, A_7 , resulting in a probability error of 0.19. In both cases the difference between the two error is below 0.1, which is the value specified for the accuracy probability threshold, γ .

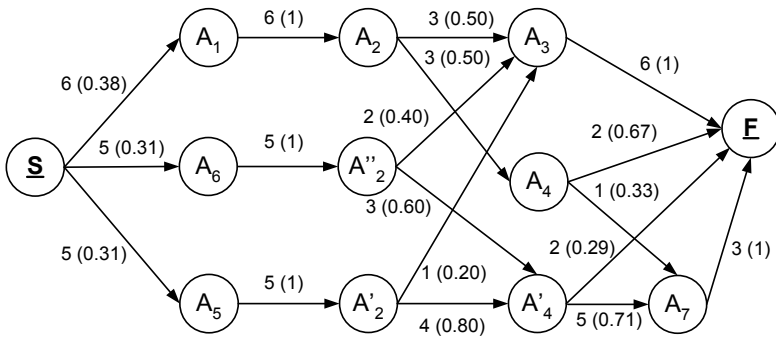


Fig. 3. The third-order model obtained when applying the dynamic clustering method with $\gamma = 0.1$ to the model given in Figure 2

Alternatively, the initial probability of a state can be estimated as $w_i / \sum_j w_j$ and every state has a link from S . In Figure 3 the error is a total of 54 page views, and, for example, $p_{S,2} = 6/54 = 0.11$ and $p_{S,4'} = 7/54 = 0.13$. For the states A_2, A_4, A_7 the probability error is given by the sum of the probabilities of the states A_2, A_4, A_7 , the states A_2', A_4', A_7 and the states A_2'', A_4', A_7 , which is 0.12. In the transition, shown in Figure 1, we have a total of 54 3-gamma counts (including 3-gamma counts starting with S and ending with F) and the count of A_2, A_4, A_7 is 6, therefore, the error is $6/54 = 0.11$. For the states A_1, A_2, A_3 the model gives 0.05, while the transition analysis gives 0.05. Both cases are accurate with respect to $\gamma = 0.1$.

The pseudo-code description of the algorithm, which implements the method, is now given. We let n be the order with which to evaluate the model, $HPG_{(n-1)}$ be the previous order model, and $(n + 1)$ -gamma be the n -gamma count of i e $n + 1$.

```

Algorithm ( $HPG_{(n-1)}, n, \gamma, (n+1)$ -grams)
begin:
  for each state  $x$ 
    induce in-paths of length  $n-1$  to  $x$ 
    for each in-path  $\vec{l}$ 
      for each out-link  $i$  from  $x$ 
        estimate  $p_{\vec{l},xi}$  and store in  $v_{\vec{l}}$ 
        if ( $|p_{\vec{l},xi} - p_{x,i}| > \gamma$ ) the state needs to be cloned
      end for
    end for
  if state needs to be cloned
    apply  $k$ -means to collection of vectors  $v_{\vec{l}}$ 
    for depth  $d = (n-1)$  to  $d = 0$ 
      for each state  $y$  at depth  $d$ 
         $c_y =$  num. distinct path prefixes assigned to different clusters
        create  $c_y - 1$  clones of state  $y$ 
        for each in-path  $\vec{l}$  to  $x$ 
          if ( $\vec{l}_{[d]} = y$  and  $\vec{l}_c > 1$ ) redirect link to corresponding clone
        end for
        for state  $y_i$  with  $i = \{1, \dots, c_y\}$ 
          for each in-link  $t$  to  $y_i$ 
            for each out-link  $r$  from  $y_i$ 
               $w_{t,y_i} = w_{t,y_i} + w_{t,y_i,r}$ ,  $w_{y_i,r} = w_{y_i,r} + w_{t,y_i,r}$ 
            end for
          end for
          remove out-links from  $y_i$  such that  $w_{y_i,r} = 0$ 
        end for
        update ngram counts to take into account clones
      end for
      update in-paths with clone references
    end for
  end if
end for
end.

```

4 Experimental Evaluation

For the experimental evaluation we analysed the real world data sets. By using data from different sources we aim to analyse the characteristics of the model in a wide enough variety of scenarios. Our previous experience has shown that individual locations and data sets have different characteristics of real world data, and therefore looking at several data sets is necessary.

The real data set (CS) information was made available by the authors of [15] and released on the website of the agency in 2002. The size was

cookie based, page caching was prohibited and data was made available with the session identified. We list the data sets in order below to enhance analysis in a wide variety of scenarios. The second data set (MM) was obtained from the authors of [11] and corresponds to one month of usage from the Magic Machine site (machines.hyperreal.org) in 1999. The data was organized in sessions and caching was disabled during collection. We list the data sets in order below, each corresponding to a week of usage. The third data set (LTM) represents forty days of usage from the London Transport Museum website in 2003 (www.ltmuseum.co.uk). The data was obtained in a raw format. We used .gif and .jpg extensions, and sessions with an empty session code. Sessions were defined as consecutive sessions for a given IP address within a 30 minute window and a maximum session length of 100 sessions was used. We list the data sets in order below, each corresponding to ten days of usage data.

Table 2 gives the summary characteristics for each data set; sessions for the data sets, γ , give the number of distinct pages visited, $\%1$ and $\% \leq 2$ indicate, respectively, the percentage of pages with just one visit and with two or less visits. Also, γ give the average number of off-line sessions, σ the standard deviation, γ the average number of in-line sessions and σ the standard deviation. Finally, γ give the number of sessions, σ the average session length, σ the standard deviation, and γ the overall number of sessions. The variability on the number of pages induced by the model for a given website can be explained by the number of pages with less than one visit. Also, when the number of pages with few visits increases the average number of off-line and in-line decreases. The average session length is highly variable because of the standard deviation. However, the MM data has a higher variability on the session length.

Table 2. Summary characteristics of the real data sets

<i>ds</i>	<i>pg</i>	$\%1v$	$\% \leq 2v$	<i>aOL</i>	<i>sOL</i>	<i>aIL</i>	<i>sIL</i>	<i>ses</i>	<i>aSes</i>	<i>sSes</i>	<i>Req</i>
LTM ₁	2998	0.62	0.68	4.5	9.6	4.4	11.6	9743	7.6	13.5	74441
LTM ₂	1648	0.19	0.27	8.4	13.8	8.3	16.6	11070	7.4	13.2	82256
LTM ₃	1610	0.27	0.37	7.8	12.8	7.7	15.0	9116	7.7	13.1	70558
LTM ₄	1586	0.24	0.34	7.8	13.3	7.7	15.9	9965	7.8	13.4	78179
MM ₁	8715	0.30	0.45	4.7	12.4	4.6	14.1	14734	6.4	37.8	94989
MM ₂	5356	0.32	0.44	6.0	18.9	5.9	20.7	14770	6.1	14.7	90682
MM ₃	5101	0.26	0.38	6.0	15.6	5.9	17.7	10924	6.7	35.2	73378
MM ₄	6740	0.35	0.49	5.1	18.5	4.9	19.8	14080	6.3	23.8	88053
CS ₁	3128	0.52	0.67	3.4	10.1	3.1	10.4	7000	4.8	6.5	33854
CS ₂	3946	0.59	0.74	2.8	9.3	2.6	9.9	7000	5.0	8.4	34897
CS ₃	5028	0.62	0.76	2.8	9.4	2.6	11.6	6950	5.5	12.8	38236

The left-hand side of Figure 4, however, for the hidden session data set, the variation of the model is with respect to $\gamma = 0$. (The data sets for each scenario reveal almost identical behavior). For the MM₁ data set a large percentage of page cloning is observed for second and third-order probabilities

which indicate that the decision can differ between high-order, ob-
 abili e and the corresponding high-order, obabili e, and that the MM1 is
 only a little better than the other when deciding which line to follow. The CS data set
 shows a lower increase in the model size, and the model can be seen to reach
 close to full accuracy with the second-order, high-order, obabili e. Finally, the LTM
 data set shows an increase in the number of states for the even high-order,
 obabili e, meaning that the choice of which line to follow is clearly influenced
 by a relatively long sequence of previously visited web pages.

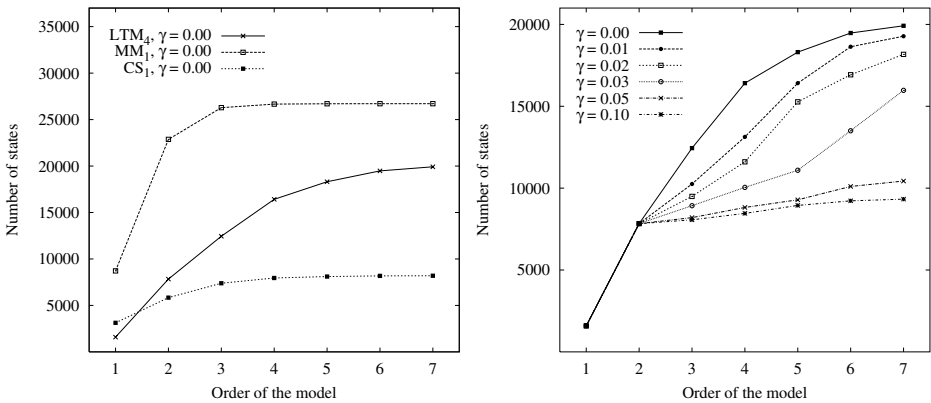


Fig. 4. The increase in model size with the model’s order for $\gamma = 0$ and the increase in size for several values of the probability threshold, γ , with the LTM₄ data set

The right-hand side of Figure 4, shows the effect of γ , on the model size for the LTM₄ data set and the left-hand side of Figure 5, shows the effect for the CS₁ data set. In both cases it can be seen that by varying the value of γ it is possible to control the increase on a model’s number of states. For both data sets, the difference in the number of states is not evident for second-order, model. For high and high-order, model it is possible to detect the number of states induced by the method by allowing some tolerance on decreasing the conditional obabili e (by setting $\gamma > 0$). Setting γ to a value greater than 0.1, results in almost no cloning for high-order, model.

Figure 6 shows a comparison on the number of clones created for the LTM₄ and CS₁ data sets, with $\gamma = 0.02$. The average number of clones created (avg) is higher for the LTM₄ data set than for the CS₁ data set, as expected by inspecting the left side of Figure 4. The standard deviation (std) indicates a substantial variability in the number of clones created, a fact highlighted by the axis number of clones (max) and the indicated percentile. For the LTM₄ data set 50% of the pages were never cloned and 75% have at least one clone for the even high-order. In the CS₁ data set 75% of the pages were never cloned and 90% of the pages have at least one even clone for the even high-order.

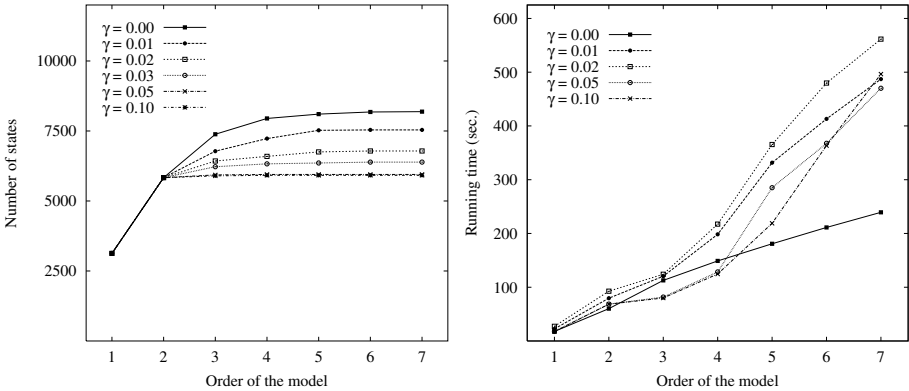


Fig. 5. The increase in model size with the model’s order for several values of the probability threshold, γ , for the CS₁ data set and variation of the running time with the model’s order for several values of the probability threshold for the LTM₄ data set

o, de . The e, e l. hel o. o i va e o , in e e in he dyna ic odel, ince while fo . o e. a e he lin choice of he co, e onding age de end on he naviga ion hi o, y, fo o he . a e he lin choice i co l e ly inde enden of he naviga ion hi o, y.

	LTM ₄ $\gamma = 0.02$					
	order					
avg	3.94	4.99	6.32	8.63	9.67	10.46
stdev	10.86	15.29	24.85	40.88	47.20	50.69
max	205	307	683	989	1193	1265
75%	4.00	5.00	5.00	6.00	6.00	6.00
85%	10.00	11.25	13.00	14.00	16.25	18.00

	CS ₁ $\gamma = 0.02$					
	order					
avg	0.87	1.06	1.11	1.16	1.17	1.17
stdev	5.40	7.06	7.3	7.92	7.96	7.96
max	138	175	180	208	208	208
75%	0.00	0.00	0.00	0.00	0.00	0.00
95%	4.00	5.00	5.00	5.00	5.00	5.00

Fig. 6. Statistics on the number of clones per state with the model’s order for the LTM₄ and CS₁ data set with $\gamma = 0.02$

The igh -hand ide of Fig , e 5, and he lef -hand ide of Fig , e 7, show o , analy i of he , nning i e of he algo, i h fo, wo , e , e en a i ve da . e . We no e ha , while , og a , ing he . e hod, we did no a e a, ic la, ca, e, e ga, ding he i l e en a ion e ciency. The . e hod i clo e o linea, i e fo, $\gamma = 0$, ince in . ch ca e no cl . e ing i needed. Fo, $\gamma > 0$ he k - ean . e hod i a l ied and we le k inc, ea e n il a ol ion which . ee . he h, e hold c i, e i a i o b ained. Fo, he , e o, ed ex e, i en . , we le k o va, y acco, ding o he ex , e ion $k = ceiling(1.5k)$ in o, de, o o b a in a low inc, e en on i . val e in he , . . age and a la, ge, inc, ea e of he k val e in he . b e en . age . Finally, he igh -hand ide of Fig , e 7, show , fo, he LTM₁ da a . e , he inc, ea e in n . be, of . a e wi h he . odel’ o, de, fo, h, ee . e hod . ed

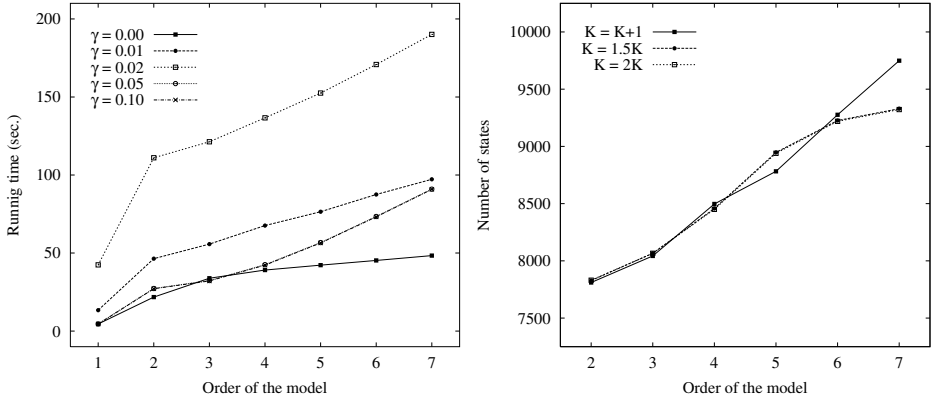


Fig. 7. The running time with the model’s order for several values of γ for the CS₁ data set and the increase in model size with the model’s order for three methods to set the number of clusters (k) for the LTM₁ data set

o increase k , in the k -mean method. The result shows how a low-order, dense network is a more informative one for the method, and for a dense, highly branched network, the method overestimates the number of clusters. This is beneficial for a fast, simple method (which is not shown in the plot).

5 Concluding Remarks

We have proposed a generalization of the HPG model by using a simple cloning operation that is able to accurately model high-order, conditional navigability. The resulting dynamic high-order Markov model is rich in the navigability of the online flow for a given parameter, effective n -order, conditional navigability of the path of the user. Thus, the model is able to capture a variable length history of usage, where different history lengths are needed to accurately model user navigation. In addition, the method allows for a navigability held together with a clustering technique that enables the control of the number of additional parameters induced by the method at the cost of accuracy. Finally, the model maintains the fundamental properties of the HPG model, [1], providing a simple, fast, and efficient algorithm for online navigation analysis, allowing for an accurate view of user navigation.

We evaluated our experiments with the real world data sets. For the results, we can conclude that, for the web site user navigation with only a short history of the usage, previously visited (for example, the MM site) but in the site, the user holds a longer history in their usage (for example, the LTM site). The result also suggests that, in a given site, different usage patterns exist in a notion of history in order to understand and the possible options they have when deciding on which link to click on. This is a good point in the proposed

dynamic model has a model each time with the selected history. The self-indication clustering method in engineering, for example, where the number of active nodes for high order, for $\gamma = 0$, becomes manageable.

In the hope we can conduct a dynamic analysis of the indicated by different order probabilities. We also plan to find a practical comparison of between order probabilities aimed at determining if the efficient practical evidence has the additional model complexity in solving a higher order, just like the corresponding increase in the algorithm's complexity. I would also be interested to be able to see the number of clusters necessary to achieve the selected accuracy in order to proceed with the method. Finally, a comparative study with the established model is planned.

References

1. J. Borges. *A data mining model to capture user web navigation patterns*. PhD thesis, University College London, London University, 2000.
2. J. Borges and M. Levene. Data mining of user navigation patterns. In B. Masand and M. Spliliopoulou, editors, *Web Usage Analysis and User Profiling*, Lecture Notes in Artificial Intelligence (LNAI 1836), pages 92–111. Springer Verlag, 2000.
3. J. Borges and M. Levene. An average linear time algorithm for web usage mining. *Int. Jou. of Information Technology and Decision Making*, 3(2):307–319, June 2004.
4. J. Borges and M. Levene. A dynamic clustering-based markov model for web usage mining. cs.IR/0406032, 2004.
5. Eugene Charniak. *Statistical Language Learning*. The MIT Press, 1996.
6. X. Chen and X. Zhang. A popularity-based prediction model for web prefetching. *Computer*, 36(3):63–70, March 2003.
7. M. Deshpande and G. Karypis. Selective markov models for predicting web page accesses. *ACM Transactions on Internet Technology*, 4(2):163–184, 2004.
8. X. Dongshan and S. Junyi. A new markov model for web access prediction. *Computing in Science and Engineering*, 4(6):34–39, November/December 2002.
9. M. Levene and G. Loizou. Computing the entropy of user navigation in the web. *Int. Journal of Information Technology and Decision Making*, 2:459–476, 2003.
10. B. Mobasher. Web usage mining and personalization. In Munindar P. Singh, editor, *Practical Handbook of Internet Computing*. Chapman Hall & CRC Press, 2004.
11. Mike Perkowitz and Oren Etzioni. Towards adaptive web sites: conceptual framework and case study. *Artificial Intelligence*, 118(2000):245–275, 2000.
12. J. Pitkow and P. Pirolli. Mining longest repeating subsequences to predict world wide web surfing. In *Proc. of the 2nd Usenix Symposium on Internet Technologies and Systems*, pages 139–150, Colorado, USA, October 1999.
13. Ramesh R. Sarukkai. Link prediction and path analysis using markov chains. *Computer Networks*, 33(1-6):377–386, June 2000.
14. S. Schechter, M. Krishnan, and M. D. Smith. Using path profiles to predict http requests. *Computer Networks and ISDN Systems*, 30:457–467, 1998.
15. M. Spiliopoulou, B. Mobasher, B. Berendt, and M. Nakagawa. A framework for the evaluation of session reconstruction heuristics in web usage analysis. *IN-FORMS Journal on Computing*, (15):171–190, 2003.
16. J. Zhu, J. Hong, and J. G. Hughes. Using markov models for web site link prediction. In *Proc. of the 13th ACM Conf. on Hypertext and Hypermedia*, pages 169–170, June 2002.

TREE² - Decision Trees for Tree Structured Data

Björn Bringmann and Albrecht Zimmermann

Institute of Computer Science, Machine Learning Lab,
Albert-Ludwigs-University Freiburg, Georges-Köhler-Allee 79,
79110 Freiburg, Germany
{bbringma, azimmerm}@informatik.uni-freiburg.de

Abstract. We present TREE², a new approach to *structural classification*. This integrated approach induces decision trees that test for pattern occurrence in the inner nodes. It combines state-of-the-art tree mining with sophisticated pruning techniques to find the most discriminative pattern in each node. In contrast to existing methods, TREE² uses no heuristics and only a single, statistically well founded parameter has to be chosen by the user. The experiments show that TREE² classifiers achieve good accuracies while the induced models are smaller than those of existing approaches, facilitating better comprehensibility.

1 Introduction

Classification is one of the most important data mining tasks. Whereas traditional approaches have focused on flat representations, using feature vectors or attribute-value representations, there has recently been a lot of interest in more expressive representations, such as sequences, trees and graphs [1,2,3,4,5]. Motivations for this interest include drug design, since molecules can be represented as graphs or sequences. Classification of such data paves the way towards drug design on the screen instead of extensive experiments in the lab. Regarding documents, XML, essentially a tree-structured representation, is becoming ever more popular. Classification in this context allows for more efficient dealing with huge amounts of electronic documents.

Existing approaches to classifying structured data (such as trees and graphs) can be categorized into various categories. They differ largely in the way they derive structural features for discriminating between examples belonging to the different classes.

A first category can be described as a pure *propositionalization* approach. The propositionalization approach typically generates a very large number of features and uses an attribute-value learner to build a classifier. The resulting classifiers are often hard to understand due to the large number of features used which are possibly also combined in a non-trivial way (e.g. in a SVM).

A second class of systems can be described as the *rule-based* approach, e.g. Zaki [4]. Even though the resulting rules often yield high predictive accuracy, the number of generated rules typically explodes, making the resulting classifier difficult to understand.

Both the association rule and propositionalization approaches consider feature generation and classification in two independent steps. There are also approaches that form a third category of systems that integrates feature construction with classification. This category includes inductive logic programming systems, such as FOIL [6] and PROGOL [7], as well as the DT-GBI approach of Motoda et al. [5]. For those approaches to be computationally feasible they have to perform heuristic search, possibly generating non-optimal features.

All techniques mentioned above share the need to specify a number of user-defined parameters, which is often non-trivial.

In this work we present a different approach called TREE². It is motivated by recent results on finding correlated patterns, allowing to find the k best features according to a convex optimization criterion such as χ^2 or \log -likelihood ratio [8]. Rather than generating a large number of features or searching for good features in a heuristic manner, TREE² searches for the best features to be incorporated in a decision tree by employing a branch-and-bound search, pruning w.r.t. the best pattern seen so far. As in DT-GBI, a decision tree is induced but at each node, the k best feature is computed. There are several advantages: TREE² is an exact algorithm, has stronger guarantees than GBI, only one parameter has to be set (the significance level), and the resulting classifiers are far smaller and easier to understand than those of the propositionalization and association rule approaches.

The paper is organized as follows: in Section 2 we describe earlier work on the topic and relate it to our approach; in Section 3, we discuss technical aspects of our method and outline our algorithm; in Section 4, the experimental evaluation is explained and its results discussed. We conclude in Section 5 and point to future work directions.

2 Related Work

Feature selection has been done with different techniques. Firstly, there are several propositionalization approaches, e.g. [2] and [3]. While details may differ, the basic mechanism in these approaches is to first mine all patterns that are unexpected according to some measure (typically frequency). Once those patterns have been found, instances are transformed into bitstrings, denoting occurrence of each pattern. Classifiers are trained using this bitstring representation. While these approaches can show excellent performance and have access to the whole spectrum of machine learning techniques there are possible problems. Obviously the decision which patterns to consider special, e.g. by fixing a minimum frequency, will have an effect on the quality of the model. The resulting feature set will probably be very large, forcing pruning of some kind. Finally, interpretation of the resulting model is not easy, especially if the classifier is non-symbolic, e.g. a SVM.

A second group of approaches is similar to the propositionalization approach [9]. Again, outstanding patterns are mined but each of them has to associate with the class value. Zaki et al.'s XRULES classifier is of this variety. Each

pattern is then considered as a rule predicting its class. Usually, the resulting rule set has to be post-processed and/or a conflict resolution technique employed. As in the propositionalization techniques, the choice of constraints under which to mine is not straight-forward and choosing the resolution technique can strongly influence performance, as has been shown e.g. in [10,11]. Additionally, the resulting classifier often consists of thousands of rules, making interpretation by the user again difficult.

Finally, there exist integrated techniques that do not mine \mathcal{L} patterns, but construct features during building the classifier. Since structural data can be represented in predicate logic, techniques such as FOIL [6] and PROGOL [7] are capable of doing that. While ILP approaches are elegant and powerful, working on large datasets can be too computationally expensive. An approach such as DT-GBI [5], on the other hand, constructs the features it uses for the tests of the induced decision tree by doing graph-mining. What is common to these approaches is that feature induction is usually done in a heuristic way, often by greedy maximization of a correlation measure during beam search. Responsibility of deciding the parameters governing this search is placed upon the user. For instance, in FOIL decisions have to be made on the beam size and the maximum number of literals that are allowed in the rule body. Similarly, DT-GBI requires the user to specify beam size, the maximum number of specializations in each node, and possibly a minimum frequency that should not be violated. As Motoda shows in his work [5], finding the right value for the beam size and the maximum number of specializations requires essentially a meta-search in the space of possible classifiers.

In contrast, the only parameter to be specified for TREE² is the cut-off value for growing the decision tree. By basing this value on the p-values for the χ^2 -distribution, the user has a well-founded guide-line for choosing this value.

While all the above techniques focus on directly using structural information for classification purposes, a different approach is exemplified by [12]. Instead of explicitly representing the structures used, kernels are employed that quantify similarities between entities. While the resulting classifiers are very accurate, the use of e.g. a graph kernel together with an SVM make analyzing the model difficult.

3 Methodology

In this section we explain the pattern matching notion used by the TREE² approach, discuss upper bound calculation, the main component of the principled search for the most discriminating pattern, and formulate the algorithm itself.

3.1 Matching Embedded Trees

Several representations for structured data such as graphs, trees and sequences exist. In this paper we will focus on tree structured data, like XML, only. Thus, we need a notion for matching tree structured data.

A rooted k -tree t is a set of k nodes V_t where each $v \in V_t$, except one called root, has a parent denoted $\pi(v) \in V_t$. We use $\lambda(v)$ to denote the label of a node and an operator \prec to denote the order from left to right among the children of a node. The transitive closure of π will be denoted π^* . Let \mathcal{L} be a formal language composed of all labeled, ordered, rooted trees and $\mathcal{D} \subset \mathcal{L}$ a database. To count trees $t \in \mathcal{D}$ containing a pattern p we define a function $d_t : \mathcal{L} \rightarrow \{0, 1\}$ to be 1 iff p matches the tree t and 0 otherwise.

Several notions of tree matching exist. As in Zaki's work [4] we used a notion called *embedded tree*, which is defined as follows:

Definition 1. A tree t is embedded in t' if there exists a mapping $\varphi : V_t \rightarrow V_{t'}$ such that $\forall u, v \in V_t : \lambda(u) = \lambda(\varphi(u)) \wedge u \prec v \Leftrightarrow \varphi(u) \prec \varphi(v) \wedge \pi^*(u) = v \Leftrightarrow \pi^*(\varphi(u)) = \varphi(v)$

An example of an embedded tree is given in Figure 1.

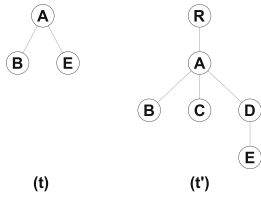


Fig. 1. The tree t is embedded in t'

	c_1	c_2	
T	y_T	$x_T - y_T$	x_T
$\neg T$	$m - y_T$	$n - m - (x_T - y_T)$	$n - x_T$
	m	$n - m$	n

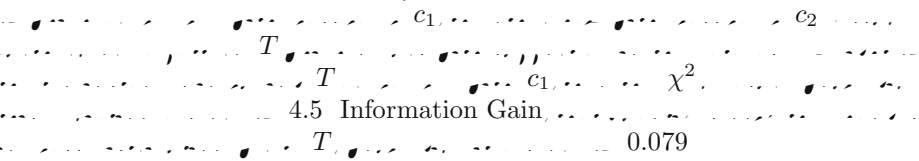
Fig. 2. A Contingency Table

We use *embedded tree* to compare our approach with Zaki's technique. This notion is more flexible than simple subtrees and the mining process is still efficient. In general, other matching notions (see [1]) and even different representations could be used with our technique. This includes not only other notions of matching trees, but also graphs, sequences etc., since the general principles of our approach apply to all domains.

3.2 Correlation Measures

Popular approaches to finding relevant patterns in the data are based on the support-confidence framework, mining frequent patterns, in the hope of capturing statistically significant phenomena, with high predictive power. This framework has some problems though, namely the difficulty of choosing a "good" support and the fact that confidence tends to reward patterns occurring together with the majority class. To alleviate these problems, we use correlation measures for selecting discriminative patterns. A correlation measure compares the expected frequency of the joint occurrence of a pattern and a certain class value to the observed frequency. If the resulting value is larger than a certain threshold then the deviation from the independence assumption is considered statistically significant enough to assume a relationship between pattern and class label.

Example 1.



We organize the observed frequencies of a tree pattern T in a contingency table, cf. Figure 2, with x_T denoting the total number of occurrences in the dataset and y_T the occurrences in the subset corresponding to the first class. Since the two variables are sufficient for calculating the value of a correlation measure on this table, we will view these measures as real-valued functions $\sigma : \mathbb{N}^2 \mapsto \mathbb{R}$ for the remainder of this paper.

While calculating the correlation value of a given pattern is relatively simple, directed search towards better solutions is somewhat more difficult since correlation measures have no desirable properties such as \dots . But if they are convex it is possible to calculate an upper bound on the score that can be achieved by specializations of the current pattern T and thus to decide whether this branch in the search tree should be followed.

3.3 Convexity and Upper Bounds

It can be proved that χ^2 and \dots are convex. For the proofs of the convexity of χ^2 and \dots we refer the reader to [8].

Convex functions take their extreme values at the points forming the convex hull of their domain D . Consider the graph of $f(x)$ in Figure 3(A). Assume the function’s domain is restricted to the interval $[k, l]$ which also makes those points the convex hull of D . Obviously, $f(k)$ and $f(l)$ are locally maximal, with $f(l)$ being the global maximum. Given the current value of the function at $f(c)$ and assuming that it is unknown whether c increases or decreases, evaluating f at k and l allows to check whether it is possible for any value of c to put the value of f over the threshold.

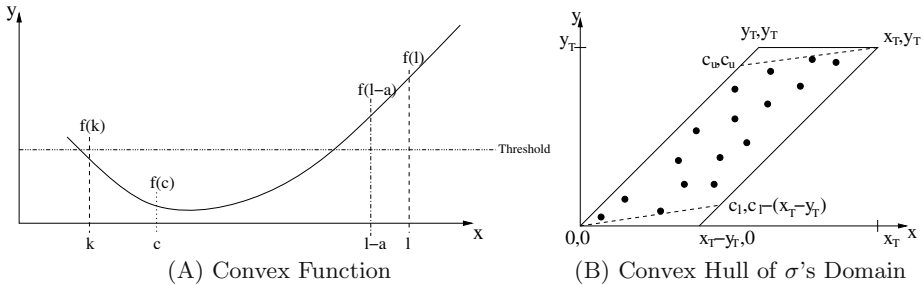


Fig. 3. Convex Function and Convex Hull of the set of possible $\langle x'_T, y'_T \rangle$

For the two-dimensional case, the extreme values are reached at the vertices of the enclosing polygon (in our case the four vertices of the parallelogram in Figure 3(B)). This parallelogram encloses all possible tuples $\langle x'_T, y'_T \rangle$ that correspond to occurrence counts of specializations of the current pattern T . The tuple $\langle 0, 0 \rangle$ corresponds to a pattern that does not occur in the dataset and therefore does not have to be considered in calculating the upper bound. $\langle x_T, y_T \rangle$ represents a valid pattern, but in the context of upper bound calculation denotes a specialization of the current pattern T that is equally good in discriminative power. Since general structures have a higher expected probability of being effective on unseen data, we prefer those and thus disregard this tuple as well. Thus the upper bound on $\sigma(T')$ is $ub_{\sigma}(T) = \max\{\sigma(y_T, y_T), \sigma(x_T - y_T, 0)\}$. For an in-depth discussion of upper bound calculation we refer the reader to [8,11].

Example 2. χ^2 , $ub_{\chi^2}(T) = \max\{9.52, 2.08\}$, $x = 10$, $y = 8$, 9.52 , $\chi^2(x_T, y_T) = 4.5$, T

While this upper bound calculation is correct for χ^2 , an additional problem w.r.t. χ^2 lies in the fact that the information provided by the score of χ^2 is not always reliable. Statistical theory says that for a contingency table with one degree of freedom, such as the one we are considering here, the expected number of occurrences has to be greater than or equal to 5 for the χ^2 score to be reliable. This means that a χ^2 -value on $\langle y_T, y_T \rangle$ or $\langle x_T - y_T, 0 \rangle$ is not necessarily reliable. Thus, upper bound calculation has to be modified to achieve reliability. Based on the size of the class and of \mathcal{D} , upper and lower bounds c_u, c_l on x'_T for which all four cells have an expected count of 5 can be calculated and the values of the tuples adjusted accordingly. Two of the new vertices are shown as $\langle c_u, c_u \rangle$ and $\langle c_l, c_l - (x_T - y_T) \rangle$.

3.4 The TREE² Algorithm

The TREE² algorithm (shown as Algorithm 1) constructs a binary decision tree in the manner of ID3 [13]. In the root node and each inner node, the occurrence of a tree pattern is tested against the instance to be classified. A resulting tree could look like the example given in Figure 4. In each node, the subtree having the best discriminative effect on the corresponding subset is found by a systematic branch-and-bound search. The mining process is shown in the subroutine ENUMERATEBESTSUBTREE. The space of possible patterns is traversed using canonical enumeration and the value of σ calculated for each candidate pattern. If this value lies above the best score seen so far, the current pattern is the most discriminating on this subset so far and the threshold is raised to its σ -value. An upper bound on the value specializations of the current pattern can achieve is calculated and pruning of the search space using this upper bound and the threshold is performed. In this way, we separate the success of the technique from user decisions about the search strategy. The only decision a user has to

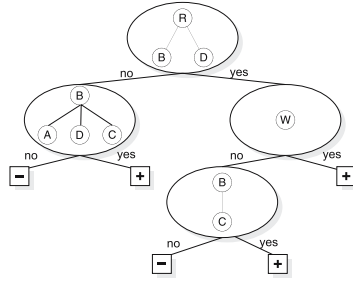


Fig. 4. A decision tree as produced by the TREE² algorithm

make is the one w.r.t. a stopping criterion for further growth of the tree. To this effect, a minimum value for the score of the correlation measure has to be specified, which can be based on statistical theory, thus giving the user a better guidance for making this decision.

Algorithm 1 The TREE² algorithm

```

TREE2( $\mathcal{D}, \sigma, \tau, \tau_{min}, DT$ )
1:  $p = \text{ENUMERATEBESTSUBTREE}(\mathbb{T}, 0, \sigma, \tau_{min}, \emptyset)$ 
2: if  $p \neq \emptyset$  then
3:   Add node including  $p$  to the DT
4:   TREE2(  $\{T \in \mathcal{D} | p \text{ embedded in } T\}, \sigma, \tau_{min}, DT$ )
5:   TREE2(  $\{T \in \mathcal{D} | p \text{ not embedded in } T\}, \sigma, \tau_{min}, DT$ )
6: return DT

ENUMERATEBESTSUBTREE( $t, \tau, \sigma, \tau_{min}, p$ )
1: for all canonical expansions  $t'$  of  $t$  do
2:   if  $\sigma(t') > \tau \wedge \sigma(t') \geq \tau_{min}$  then
3:      $p = t', \tau = \sigma(t')$ 
4:   if  $ub(t') \geq \tau$  then
5:      $p = \text{ENUMERATEBESTSUBTREE}(t', \tau, \sigma, \tau_{min}, p)$ 
6: return  $p$ 

```

TREE² has several desirable properties. Firstly, the resulting classifier is integrated in the sense that it uses patterns directly, thus circumventing the need for the user to restrict the amount of features and making the resulting classifier more understandable. Secondly, by using correlation measures for quantifying the quality of patterns, we give the user a sounder theoretical foundation on which to base the decision about which learned tests to consider significant and include in the model. Thirdly, we avoid using heuristics that force the user to decide on the values of parameters that could have a severe impact on the resulting model’s accuracy. Using principled search guarantees that TREE² finds the best discriminating pattern for each node in the decision tree w.r.t. the correlation measure used. Finally, as the experiments show, the resulting decision tree is far smaller than the rule sets produced by XRULES classifier [4], while achieving comparable accuracy, and is therefore more easily interpretable by human users.

4 Experimental Evaluation

For the experimental evaluation, we compared our approach to XRULES and a decision tree base-line approach on the XML data used in Zaki et al.'s publication [4]. Furthermore, we compared TREE² to a base-line approach using frequency mining for a SVM classifier and two PROGOL results on the regression-friendly subset of the Mutagenesis dataset.

XML Data. The XML data used in our experiments are log files from web-site visitors' sessions. They are separated into three weeks (CSLOG1, CSLOG2, and CSLOG3) and each session is classified as its producing visitor coming either from an .edu domain or from any other domain. Characteristics of the datasets are shown in Table 1. For the comparison we built decision trees with the

Table 1. Characteristics of Datasets (taken from [4])

DB	#Sessions	edu	other	%edu	%other
CSLOG1	8074	1962	6112	24.3	75.7
CSLOG2	7409	1687	5722	22.8	77.2
CSLOG12	13934	2969	10965	21.3	78.7
CSLOG3	7628	1798	5830	23.6	76.4

χ^2 distribution's significance value for 90%, 95% and 99% respectively. In each setting we used one set of data for training and another one for testing. Following Zaki's notation, CSLOG x - y means that we trained on set x and tested on set y . For the base-line approach we mined the 100 patterns having the highest discriminative effect on the data, transformed the data into bitstring instances according to the found patterns, and built decision trees using all 100 patterns in one run () and the 50 best patterns in another run () with the WEKA [14] implementation of the C4.5 [15] algorithm. We compare the accuracies of the resulting classifiers against each other as well as the complexity of the model which we measure by the number of rules used by XRULES, and by the number of leaves in the decision trees, which corresponds to the number of rules that can be derived from the trees, respectively.

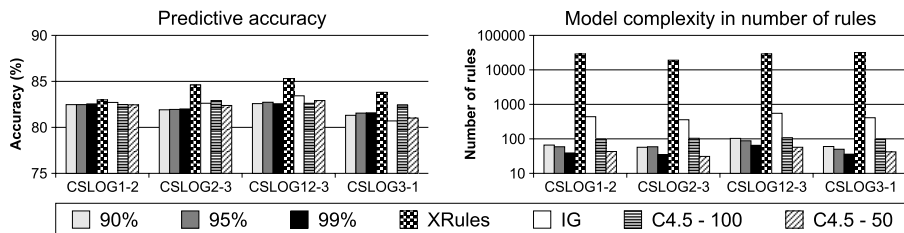


Fig. 5. Accuracies and size in rules of the different approaches

Results are summarized in Figure 5. As can be seen, the accuracies of the induced classifiers do not vary much. The only approach that significantly outperforms (by 2-3%) the other techniques on all but the CSLOG1-2 setting, is XRULES. At the same time, the size of XRULES' models is also significantly greater. While the TREE² trees induced with χ^2 have several hundred nodes and all trees induced with χ^2 (both TREE² and base-line) between 35 and 103 nodes, the smallest XRULES model consists of more than 19000 rules. Patterns tested against in the inner decision tree nodes consist of 3-7 nodes only. Since this is similar to the size of patterns used in XRULES' rules, complexity is really reduced and not just pushed inside the classifier. In comparing the other approaches, several things are noticeable. Raising the threshold from the 90% to the 95% significance level for χ^2 -induced TREE² trees does not decrease accuracy (even improving it slightly in 3 cases). Raising it further to the 99% level has no clear effect. The tree size decreases, though, on average by 7.5 nodes from the 90% to the 95% setting. Raising the significance level further to 99% decreases the tree size by 18 nodes on average.

For the base-line approach we mined patterns correlating strongly with the classes and trained a classifier on them. This approach achieves competitive results w.r.t the accuracy. The clear drawback is that deciding on the number of features to use is not straightforward. Using only 50 instead of 100 features produces all kinds of behavior. In some cases the accuracy does not change. In other cases the classifier using 50 features outperforms the one using 100 or vice versa. Also, the base-line approach using 100 patterns tends to use most of these, even if TREE² trees of similar quality are much smaller.

Finally, using χ^2 as quality criterion shows mainly one thing - that it is difficult to make an informed decision on cut-off values. The accuracies and sizes shown refer to decision trees induced with a cut-off value of 0.001. For one thing, the resulting trees grow far bigger than the χ^2 -trees. Additionally, the accuracies in comparison with the χ^2 approach vary, giving rise to one worse tree, one of equal quality and two better ones. None of the differences in accuracy is significant though. Inducing decision trees with a cut-off value of 0.01 lowers accuracy by 1.5 to 3 percentage points, with the induced trees still being larger than the χ^2 trees.

Mutagenicity Data. For this setting, we chose the regression-friendly subset of the well known Mutagenicity dataset used in [16]. We compare with the results of the ILP system PROGOL reported in [16,17] and the results of the base-line approach reported in [3]. Since the Mutagenicity dataset consists of molecules represented as graphs, a transformation from the SMILES representation into so-called fragment-trees is used that is explained following this paragraph.

The SMILES language [18] is used by computational chemists as a compact encoding of molecular structure. It is supported by many tools as OpenBabel or Daylight ([19,20]). The language contains symbols for atoms, bonds, branches, and can express cycles. Using a decomposition-algorithm by Karwath and De Raedt [21], a SMILES-String can, after some

reformatting, be decomposed into a so-called \dots . Since there is no \dots SMILES-string for a molecule, the fragment tree is not unique either. The decomposition-algorithm recursively splits the string into $\dots \{xT\}_x$ and $\dots A(B)C$. In the resulting fragment-tree the leaves contain pure cycles or linear fragments without further branches. The inner nodes of such a tree contain fragments still containing branches while the root node is the whole molecule. The edge labels denote the type of decomposition (i.e. the part of the branch or the number of the cycle). Thus, the leaves of a fragment-tree contain a lot of information decomposed into very small fragments. As in [3] we drop the edge labels and labeled all but the leaf nodes with a new, unique label. Hence, the tree-structure represents the abstract structure of the molecule with the chemical information in the leaves.

Figure 6 shows a molecule on the left-hand side which could be encoded by the SMILES-string $N - c1ccc(cc1) - O - c2ccc(cc2) - [Cl]$. This string represents the same as $N\{0cccc(cc)_0\}O\{1cccc(cc)_1\}[Cl]$. The corresponding fragment-tree is shown on the right-hand side of Figure 6.

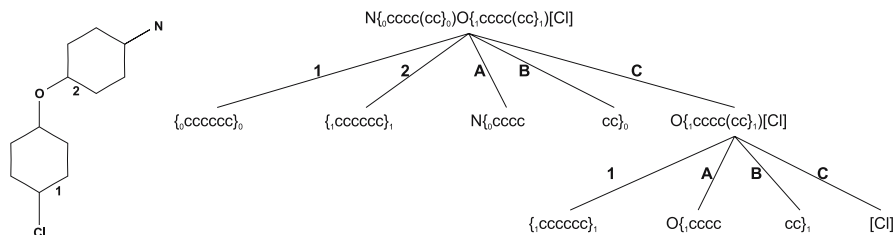


Fig. 6. A molecule with the encoding $N - c1ccc(cc1) - O - c2ccc(cc2) - [Cl]$ and the corresponding fragment-tree

\dots Predictive accuracy for each approach was estimated using ten-fold cross-validation. Reported are average accuracies and standard deviation (if known). For TREE², trees were induced at the 95% significance level for χ^2 and with a cut-off value of 0.01 for \dots . The results reported in [16] were achieved using PROGOL and working only on structural information, in [17], numerical values suggested by experts were used as well. This work reports only an average accuracy. The resulting accuracies and the size of the corresponding theories are shown in Table 2.

As can be seen, for both measures TREE² gives similar results to the purely structural PROGOL approach, with the differences being not significant. At the same time, the χ^2 induced model is far smaller than the other two. Again, the patterns tested against in the inner nodes are not overly complex (5-11 nodes). When PROGOL uses the expert-identified attributes as well, its accuracy increases. Since we do not have access to the standard deviation of these experiments, we cannot make a significance statement. Finally, the base-line approach,

Table 2. Accuracies and complexity of the models on the mutagenicity dataset

Approach	Predictive Accuracy	Average Size of the Model
TREE ² χ^2	80.26±7.14	2.3 Nodes
TREE ² <i>IG</i>	81.76±9	11.8 Nodes
PROGOL '94 [16]	80±3	9 Clauses
PROGOL '95 [17]	84	4 Clauses
FREQUENT SMILES [3]	86.70	214 Patterns

which mined all patterns frequent in one class and not exceeding a given frequency in the other class, and built a model using these features in an SVM, significantly outperforms the TREE² classifiers. On the other hand, by using a SVM, the results will hardly be interpretable for humans anymore and the amount of patterns used is larger than in the TREE² models by two orders of magnitude.

5 Conclusion and Future Work

We presented TREE², an integrated approach to structural classification. The algorithm builds a decision tree for tree structured data that tests for pattern occurrence in the inner nodes. Using an optimal branch-and-bound search, made possible by effective pruning, TREE² finds the most discriminative pattern for each subset of the data considered. This allows the user to abstract the success of the classifier from decisions about the search process, unlike in existing approaches that include heuristics. Basing the stopping criterion for growing the decision tree on statistically well founded measures rather than arbitrary thresholds whose meaning is somewhat ambiguous gives the user better guidance for selecting this parameter. It also alleviates the main problem of the support-confidence framework, namely the generation of very large rule sets that are incomprehensible to the user and possibly include uninformative rules w.r.t. classification.

As the experiments show, TREE² classifiers are effective while being less complex than existing approaches. While using χ^2 for assessing the quality of discriminative patterns, raising or lowering the significance threshold affects the induced trees in an expected manner. In contrast, using support-confidence is more difficult, since selecting the cut-off value has no statistical foundations. While base-line approaches, that separate feature generation and classifier construction, achieve very good results, it is not entirely clear how to justify the selected the number of features mined. Furthermore, there exists a gap in interpretability since the classifier used might combine the mined features in a way that is not easily accessible to the user.

So far, we have restricted ourselves to a single representation, a certain type of classifier, and two measures. Future work will include evaluating other correlation measures and applying our approach to different

representations. Finally, the success of using effective conflict resolution strategies in the XRULES classifier suggests the expansion our approach to ensemble classifiers.

Acknowledgments. We would like to thank Mohammed J. Zaki for providing the datasets and the XRULES algorithm. Furthermore, we would like to thank Andreas Karwath, Kristian Kersting, Robert Egginton and Luc De Raedt for interesting discussions and comments to our work.

References

1. Kilpeläinen, P.: Tree Matching Problems with Applications to Structured Text Databases. PhD thesis, University of Helsinki (1992)
2. Kramer, S., Raedt, L.D., Helma, C.: Molecular feature mining in HIV data. In Provost, F., Srikant, R., eds.: Proc. KDD-01, New York, ACM Press (2001) 136–143
3. Bringmann, B., Karwath, A.: Frequent SMILES. In: Lernen, Wissensentdeckung und Adaptivität, Workshop GI Fachgruppe Maschinelles Lernen, LWA. (2004)
4. Zaki, M.J., Aggarwal, C.C.: XRules: an effective structural classifier for XML data. In Getoor, L., Senator, T.E., Domingos, P., Faloutsos, C., eds.: KDD, Washington, DC, USA, ACM (2003) 316–325
5. Geamsakul, W., Matsuda, T., Yoshida, T., Motoda, H., Washio, T.: Performance evaluation of decision tree graph-based induction. In Grieser, G., Tanaka, Y., Yamamoto, A., eds.: Discovery Science, Sapporo, Japan, Springer (2003) 128–140
6. Quinlan, J.R.: Learning logical definitions from relations. *Machine Learning* **5** (1990) 239–266
7. Muggleton, S.: Inverse entailment and PROGOL. *New Generation Computing* **13** (1995) 245–286
8. Morishita, S., Sese, J.: Traversing itemset lattices with statistical metric pruning. In: Proceedings of the Nineteenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, Dallas, Texas, USA, ACM (2000) 226–236
9. Liu, B., Hsu, W., Ma, Y.: Integrating classification and association rule mining. In Agrawal, R., Stolorz, P.E., Piatetsky-Shapiro, G., eds.: KDD, New York City, New York, USA, AAAI Press (1998) 80–86
10. Mutter, S., Hall, M., Frank, F.: Using classification to evaluate the output of confidence-based association rule mining. In Webb, G.I., Yu, X., eds.: Australian Conference on Artificial Intelligence, Cairns, Australia, Springer (2004) 538–549
11. Zimmermann, A., De Raedt, L.: Corclass: Correlated association rule mining for classification. [22] 60–72
12. Gärtner, T., Lloyd, J.W., Flach, P.A.: Kernels and distances for structured data. *Machine Learning* **57** (2004)
13. Quinlan, J.R.: Induction of decision trees. *Machine Learning* **1** (1986) 81–106
14. Frank, E., Hall, M., Trigg, L.E., Holmes, G., Witten, I.H.: Data mining in bioinformatics using weka. *Bioinformatics* **20** (2004) 2479–2481
15. Quinlan, J.R.: C4.5: Programs for Machine Learning. Morgan Kaufmann (1993)
16. Srinivasan, A., Muggleton, S., King, R., Sternberg, M.: Mutagenesis: ILP experiments in a non-determinate biological domain. In Wrobel, S., ed.: Proceedings of the 4th International Workshop on Inductive Logic Programming. Volume 237., Gesellschaft für Mathematik und Datenverarbeitung MBH (1994) 217–232

17. King, R.D., Sternberg, M.J.E., Srinivasan, A.: Relating chemical activity to structure: An examination of ILP successes. *New Generation Comput.* **13** (1995) 411–433
18. Weininger, D.: SMILES, a chemical language and information system 1. Introduction and encoding rules. *J. Chem. Inf. Comput. Sci.* **28** (1988) 31–36
19. The OpenBabel Software Community: Open Babel. <http://openbabel.sourceforge.net/> (2003)
20. Daylight Chemical Information Systems, Inc. <http://www.daylight.com/> (2004)
21. Karwath, A., De Raedt, L.: Predictive graph mining. [22] 1–15
22. Suzuki, E., Arikawa, S., eds.: Discovery Science, 7th International Conference, DS 2004, Padova, Italy, October 2-5, 2004, Proceedings. In Suzuki, E., Arikawa, S., eds.: DS 2004, Padova, Italy, Springer (2004)

Agglomerative Hierarchical Clustering with Constraints: Theoretical and Empirical Results

Ian Davidson and S.S. Ravi

Department of Computer Science, University at Albany - State University of New York,
Albany, NY 12222
{davidson, ravi}@cs.albany.edu

Abstract. We explore the use of instance and cluster-level constraints with agglomerative hierarchical clustering. Though previous work has illustrated the benefits of using constraints for non-hierarchical clustering, their application to hierarchical clustering is not straight-forward for two primary reasons. First, some constraint combinations make the feasibility problem (Does there exist a single feasible solution?) **NP**-complete. Second, some constraint combinations when used with traditional agglomerative algorithms can cause the dendrogram to stop prematurely in a dead-end solution even though there exist other feasible solutions with a significantly smaller number of clusters. When constraints lead to efficiently solvable feasibility problems and standard agglomerative algorithms do not give rise to dead-end solutions, we empirically illustrate the benefits of using constraints to improve cluster purity and average distortion. Furthermore, we introduce the new γ constraint and use it in conjunction with the triangle inequality to considerably improve the efficiency of agglomerative clustering.

1 Introduction and Motivation

Hierarchical clustering algorithms are run once and create a dendrogram which is a tree structure containing a k -block set partition for each value of k between 1 and n , where n is the number of data points to cluster allowing the user to choose a particular clustering granularity. Though less popular than non-hierarchical clustering there are many domains [16] where clusters naturally form a hierarchy; that is, clusters are part of other clusters. Furthermore, the popular agglomerative algorithms are easy to implement as they just begin with each point in its own cluster and progressively join the closest clusters to reduce the number of clusters by 1 until $k = 1$. The basic agglomerative hierarchical clustering algorithm we will improve upon in this paper is shown in Figure 1. However, these added benefits come at the cost of time and space efficiency since a typical implementation with symmetrical distances requires $\Theta(mn^2)$ computations, where m is the number of attributes used to represent each instance.

In this paper we shall explore the use of instance and cluster level constraints with hierarchical clustering algorithms. We believe the use of such constraints with *hierarchical* clustering is the first though there exists work that uses spatial constraints to find specific types of clusters and avoid others [14,15]. The similarly named *constrained hierarchical clustering* [16] is actually a method of combining partitional and hierarchical clustering algorithms; the method does not incorporate apriori constraints. Recent work

Agglomerative($S = \{x_1, \dots, x_n\}$) **returns** *Dendrogram* $_k$ for $k = 1$ to $|S|$.

1. $C_i = \{x_i\}, \forall i$.
 2. **for** $k = |S|$ **down to** 1
 - $Dendrogram_k = \{C_1, \dots, C_k\}$
 - $d(i, j) = D(C_i, C_j), \forall i, j; \quad l, m = \operatorname{argmin}_{a,b} d(a, b)$.
 - $C_l = \operatorname{Join}(C_l, C_m); \quad \operatorname{Remove}(C_m)$.
- endloop**
-

Fig. 1. Standard Agglomerative Clustering

[1,2,12] in the non-hierarchical clustering literature has explored the use of instance-level constraints. The **must-link** and **cannot-link** constraints require that two instances must both be part of or not part of the same cluster respectively. They are particularly useful in situations where a large amount of unlabeled data to cluster is available along with some labeled data from which the constraints can be obtained [12]. These constraints were shown to improve cluster purity when measured against an extrinsic class label not given to the clustering algorithm [12]. The δ **constraint** requires the distance between any pair of points in two different clusters to be at least δ . For any cluster C_i with two or more points, the ϵ -**constraint** requires that for each point $x \in C_i$, there must be another point $y \in C_i$ such that the distance between x and y is at most ϵ . Our recent work [4] explored the computational complexity (difficulty) of the *feasibility* problem: **Given** a value of k , does there exist at least one clustering solution that satisfies all the constraints and has k clusters? Though it is easy to see that there is no feasible solution for the three cannot-link constraints $CL(a,b), CL(b,c), CL(a,c)$ for $k < 3$, the general feasibility problem for cannot-link constraints is **NP**-complete by a reduction from the graph coloring problem. The complexity results of that work, shown in Table 1 (2nd column), are important for data mining because when problems are shown to be intractable in the worst-case, we should avoid them or should not expect to find an exact solution efficiently.

We begin this paper by exploring the feasibility of agglomerative **hierarchical** clustering under the above four mentioned instance and cluster-level constraints. This problem is *significantly* different from the feasibility problems considered in our previous work since the value of k for hierarchical clustering is not given. We then empirically show that constraints with a modified agglomerative hierarchical algorithm can improve the quality and performance of the resultant dendrogram. To further improve performance we introduce the γ constraint which when used with the triangle inequality can yield large computation saving that we have bounded in the best and average case. Finally, we cover the interesting result of an irreducible clustering. If we are given a feasible clustering with k_{max} clusters then for certain combination of constraints joining the two closest clusters may yield a feasible but “dead-end” solution with k clusters from which no other feasible solution with less than k clusters can be obtained, even though they are known to exist. Therefore, the created dendrograms may be incomplete.

Throughout this paper $D(x, y)$ denotes the Euclidean distance between two points and $D(X, Y)$ the Euclidean distance between the centroids of two groups of instances. We note that the feasibility and irreducibility results (Sections 2 and 5) are not neces-

Table 1. Results for Feasibility Problems for a Given k (partitional clustering) and Unspecified k (hierarchical clustering)

Constraint	Given k	Unspecified k	Unspecified k - Deadends?
Must-Link	P [9,4]	P	No
Cannot-Link	NP -complete [9,4]	P	Yes
δ -constraint	P [4]	P	No
ϵ -constraint	P [4]	P	No
Must-Link and δ	P [4]	P	No
Must-Link and ϵ	NP -complete [4]	P	No
δ and ϵ	P [4]	P	No
Must-Link, Cannot-Link, δ and ϵ	NP -complete [4]	NP -complete	Yes

sarily for Euclidean distances and are hence applicable for single and complete linkage clustering while the γ -constraint to improve performance (Section 4) is applicable to any metric space.

2 Feasibility for Hierarchical Clustering

In this section, we examine the feasibility problem for several different types of constraints, that is, the problem of determining whether the given set of points can be partitioned into clusters so that all the specified constraints are satisfied.

Definition 1. Feasibility problem for Hierarchical Clustering (FHC)

Instance: A set S of nodes, the (symmetric) distance $d(x, y) \geq 0$ for each pair of nodes x and y in S and a collection C of constraints.

Question: Can S be partitioned into subsets (clusters) so that all the constraints in C are satisfied?

When the answer to the feasibility question is “yes”, the corresponding algorithm also produces a partition of S satisfying the constraints. We note that the FHC problem considered here is *significantly* different from the constrained non-hierarchical clustering problem considered in [4] and the proofs are different as well even though the end results are similar. For example in our earlier work we showed intractability results for some constraint types using a straightforward reduction from the graph coloring problem. The intractability proof used in this work involves more elaborate reductions. For the feasibility problems considered in [4], the number of clusters is in effect, another constraint. In the formulation of FHC, there are *no* constraints on the number of clusters, other than the trivial ones (i.e., the number of clusters must be at least 1 and at most $|S|$).

We shall in this section begin with the same constraints as those considered in [4]. They are: (a) Must-Link (ML) constraints, (b) Cannot-Link (CL) constraints, (c) δ constraint and (d) ϵ constraint. In later sections we shall introduce another cluster-level

constraint to improve the efficiency of the hierarchical clustering algorithms. As observed in [4], a δ constraint can be efficiently transformed into an equivalent collection of ML-constraints. Therefore, we restrict our attention to ML, CL and ϵ constraints. We show that for any *pair* of these constraint types, the corresponding feasibility problem can be solved efficiently. The simple algorithms for these feasibility problems can be used to seed an agglomerative or divisive hierarchical clustering algorithm as is the case in our experimental results. However, when all three types of constraints are specified, we show that the feasibility problem is **NP**-complete and hence finding a clustering, let alone a good clustering, is computationally intractable.

2.1 Efficient Algorithms for Certain Constraint Combinations

When the constraint set C contains only ML and CL constraints, the FHC problem can be solved in polynomial time using the following simple algorithm.

1. Form the clusters implied by the ML constraints. (This can be done by computing the transitive closure of the ML constraints as explained in [4].) Let C_1, C_2, \dots, C_p denote the resulting clusters.
2. If there is a cluster C_i ($1 \leq i \leq p$) with nodes x and y such that x and y are also involved in a CL constraint, then there is no solution to the feasibility problem; otherwise, there is a solution.

When the above algorithm indicates that there is a feasible solution to the given FHC instance, one such solution can be obtained as follows. Use the clusters produced in Step 1 along with a singleton cluster for each node that is not involved in an ML constraint. Clearly, this algorithm runs in polynomial time. We now consider the combination of CL and ϵ constraints. Note that there is always a trivial solution consisting of $|S|$ singleton clusters to the FHC problem when the constraint set involves only CL and ϵ constraints. Obviously, this trivial solution satisfies both CL and ϵ constraints, as the latter constraint only applies to clusters containing two or more instances.

The FHC problem under the combination of ML and ϵ constraints can be solved efficiently as follows. For any node x , an ϵ -neighbor of x is another node y such that $D(x, y) \leq \epsilon$. Using this definition, an algorithm for solving the feasibility problem is:

1. Construct the set $S' = \{x \in S : x \text{ does not have an } \epsilon\text{-neighbor}\}$.
2. If some node in S' is involved in an ML constraint, then there is no solution to the FHC problem; otherwise, there is a solution.

When the above algorithm indicates that there is a feasible solution, one such solution is to create a singleton cluster for each node in S' and form one additional cluster containing all the nodes in $S - S'$. It is easy to see that the resulting partition of S satisfies the ML and ϵ constraints and that the feasibility testing algorithm runs in polynomial time. The following theorem summarizes the above discussion and indicates that we can extend the basic agglomerative algorithm with these combinations of constraint types to perform efficient hierarchical clustering. However, it does not mean that we can always use traditional agglomerative clustering algorithms as the closest-cluster-join operation can yield dead-end clustering solutions as discussed in Section 5.

Theorem 1. *The FHC problem can be solved efficiently for each of the following combinations of constraint types: (a) ML and CL (b) CL and ϵ and (c) ML and ϵ . \square*

2.2 Feasibility Under ML, CL and ϵ Constraints

In this section, we show that the FHC problem is **NP**-complete when all the three constraint types are involved. This indicates that creating a dendrogram under these constraints is an intractable problem and the best we can hope for is an approximation algorithm that may **not** satisfy all constraints. The **NP**-completeness proof uses a reduction from the One-in-Three 3SAT with positive literals problem (OPL) which is known to be **NP**-complete [11]. For each instance of the OPL problem we can construct a constrained clustering problem involving ML, CL and ϵ constraints. Since complexity results are worse case, the existence of just these problems is sufficient for theorem 2.

One-in-Three 3SAT with Positive Literals (OPL)

Instance: A set $C = \{x_1, x_2, \dots, x_n\}$ of n Boolean variables, a collection $Y = \{Y_1, Y_2, \dots, Y_m\}$ of m clauses, where each clause $Y_j = \{x_{j_1}, x_{j_2}, x_{j_3}\}$ has exactly three non-negated literals.

Question: Is there an assignment of truth values to the variables in C so that exactly one literal in each clause becomes true?

Theorem 2. *The FHC problem is NP-complete when the constraint set contains ML, CL and ϵ constraints.*

The proof of the above theorem is somewhat lengthy and is omitted because of space reasons. (The proof appears in an expanded technical report version of this paper [5] that is available on-line.)

3 Using Constraints for Hierarchical Clustering: Algorithm and Empirical Results

To use constraints with hierarchical clustering we change the algorithm in Figure 1 to factor in the above discussion. As an example, a constrained hierarchical clustering algorithm with must-link and cannot-link constraints is shown in Figure 2. In this section we illustrate that constraints can improve the quality of the dendrogram. We purposefully chose a small number of constraints and believe that even more constraints will improve upon these results. We will begin by investigating must-link and cannot-link constraints using six real world UCI datasets. For each data set we clustered all instances but removed the labels from 90% of the data (S_u) and used the remaining 10% (S_l) to generate constraints. We randomly selected two instances at a time from S_l and generated must-link constraints between instances with the same class label and cannot-link constraints between instances of differing class labels. We repeated this process twenty times, each time generating 250 constraints of each type. The performance measures reported are averaged over these twenty trials. All instances with missing values were

Table 2. Average Distortion per Instance and Average Percentage Cluster Purity over Entire Dendrogram

Data Set	Distortion		Purity	
	Unconstrained	Constrained	Unconstrained	Constrained
Iris	3.2	2.7	58%	66%
Breast	8.0	7.3	53%	59%
Digit (3 vs 8)	17.1	15.2	35%	45%
Pima	9.8	8.1	61%	68%
Census	26.3	22.3	56%	61%
Sick	17.0	15.6	50%	59%

ConstrainedAgglomerative(S,ML,CL) returns *Dendrogram_i*, $i = k_{min} \dots k_{max}$

Notes: In Step 5 below, the term “mergeable clusters” is used to denote a pair of clusters whose merger does not violate any of the given CL constraints. The value of t at the end of the loop in Step 5 gives the value of k_{min} .

1. Construct the transitive closure of the ML constraints (see [4] for an algorithm) resulting in r connected components M_1, M_2, \dots, M_r .
2. If two points $\{x, y\}$ are both a CL and ML constraint then output “No Solution” and stop.
3. Let $S_1 = S - (\bigcup_{i=1}^r M_i)$. Let $k_{max} = r + |S_1|$.
4. Construct an initial feasible clustering with k_{max} clusters consisting of the r clusters M_1, \dots, M_r and a singleton cluster for each point in S_1 . Set $t = k_{max}$.
5. **while** (there exists a pair of mergeable clusters) **do**
 - (a) Select a pair of clusters C_l and C_m according to the specified distance criterion.
 - (b) Merge C_l into C_m and remove C_l . (The result is *Dendrogram_{t-1}*.)
 - (c) $t = t - 1$.

endwhile

Fig. 2. Agglomerative Clustering with ML and CL Constraints

removed as hierarchical clustering algorithms do not easily handle such instances. Furthermore, all non-continuous columns were removed as there is no standard distance measure for discrete columns.

Table 2 illustrates the quality improvement that the must-link and cannot-link constraints provide. Note that we compare the dendrograms for k values between k_{min} and k_{max} . For each corresponding level in the unconstrained and constrained dendrogram we measure the average distortion ($1/n * \sum_{i=1}^n D(x_i - C_{f(x_i)})$), where $f(x_i)$ returns the index of the closest cluster to x_i and present the average over all levels. It is important to note that we are not claiming that agglomerative clustering has distortion as an objective function, rather that it is a good measure of cluster quality. We see that the distortion improvement is typically of the order of 15%. We also see that the average percentage purity of the clustering solution as measured by the class label purity improves. The cluster purity is measured against the extrinsic class labels. We believe these improvement are due to the following. When many pairs of clusters have simi-

lar short distances, the must-link constraints guide the algorithm to a better join. This type of improvement occurs at the bottom of the dendrogram. Conversely, towards the top of the dendrogram the cannot-link constraints rule out ill-advised joins. However, this preliminary explanation requires further investigation which we intend to address in the future. In particular, a study of the most informative constraints for hierarchical clustering remains an open question, though promising preliminary work for the area of non-hierarchical clustering exists [2].

We next use the cluster-level δ constraint with an arbitrary value to illustrate the great computational savings that such constraints offer. Our earlier work [4] explored ϵ and δ constraints to provide background knowledge towards the “type” of clusters we wish to find. In that paper we explored their use with the Aibo robot to find objects in images that were more than 1 foot apart as the Aibo can only navigate between such objects. For these UCI data sets no such background knowledge exists and how to set these constraint values for non-spatial data remains an active research area. Hence we test these constraints with arbitrary values. We set δ equal to 10 times the average distance between a pair of points. Such a constraint will generate hundreds even thousands of must-link constraints that can greatly influence the clustering results and algorithm efficiency as shown in Table 3. We see that the minimum improvement was 50% (for Census) and nearly 80% for Pima. This improvement is due to the constraints effectively creating a pruned dendrogram by making $k_{max} \ll n$.

Table 3. The Rounded Mean Number of Pair-wise Distance Calculations for an Unconstrained and Constrained Clustering using the δ constraint

Data Set	Unconstrained	Constrained
Iris	22,201	3,275
Breast	487,204	59,726
Digit (3 vs 8)	3,996,001	990,118
Pima	588,289	61,381
Census	2,347,305,601	563,034,601
Sick	793,881	159,801

4 Using the γ Constraint to Improve Performance

In this section we introduce a new constraint, the γ constraint and illustrate how the triangle inequality can be used to further improve the run-time performance of agglomerative hierarchical clustering. Though this improvement does not affect the worst-case analysis, we can perform a best case analysis and an expected performance improvement using the Markov inequality. Future work will investigate if tighter bounds can be found. There exists other work involving the triangle inequality but not constraints for non-hierarchical clustering [6] as well as for hierarchical clustering [10].

Definition 2. (*The γ Constraint For Hierarchical Clustering*) Two clusters whose geometric centroids are separated by a distance greater than γ cannot be joined.

IntelligentDistance (γ , $C = \{C_1, \dots, C_k\}$)

returns $d(i, j) \forall i, j$.

1. **for** $i = 2$ **to** $n - 1$ $d_{1,i} = D(C_1, C_i)$ **endloop**
 2. **for** $i = 2$ **to** $n - 1$
 - for** $j = i + 1$ **to** $n - 1$ $\hat{d}_{i,j} = |d_{1,i} - d_{1,j}|$
 - if** $\hat{d}_{i,j} > \gamma$ **then** $d_{i,j} = \gamma + 1$; *do not join* **else** $d_{i,j} = D(x_i, x_j)$
 - endloop**
 - endloop**
 3. **return** $d_{i,j}, \forall i, j$.
-

Fig. 3. Function for Calculating Distances Using the γ Constraint and the Triangle Inequality

The γ constraint allows us to specify how geometrically well separated the clusters should be. Recall that the triangle inequality for three points a, b, c refers to the expression $|D(a, b) - D(b, c)| \leq D(a, c) \leq D(a, b) + D(c, b)$ where D is the Euclidean distance function or any other metric function. We can improve the efficiency of the hierarchical clustering algorithm by making use of the lower bound in the triangle inequality and the γ constraint. Let a, b, c now be cluster centroids and we wish to determine the closest two centroids to join. If we have already computed $D(a, b)$ and $D(b, c)$ and the value $|D(a, b) - D(b, c)|$ exceeds γ , then we need not compute the distance between a and c as the lower bound on $D(a, c)$ already exceeds γ and hence a and c cannot be joined. Formally the function to calculate distances using geometric reasoning at a particular dendrogram level is shown in Figure 3. Central to the approach is that the distance between a central point (c) (in this case the first) and every other point is calculated. Therefore, when bounding the distance between two instances (a, b) we effectively calculate a triangle with two edges with known lengths incident on c and thereby lower bound the distance between a and b . How to select the best central point and the use of multiple central points remains future important research.

If the triangle inequality bound exceeds γ , then we save making m floating point power calculations if the data points are in m dimensional space. As mentioned earlier we have no reason to believe that there will be at least one situation where the triangle inequality saves computation in *all problem instances*; hence in the worst case, there is no performance improvement. But in practice it is expected to occur and hence we can explore the best and expected case results.

4.1 Best Case Analysis for Using the γ Constraint

Consider the n points to cluster $\{x_1, \dots, x_n\}$. The first iteration of the agglomerative hierarchical clustering algorithm using symmetrical distances is to compute the distance between each point and every other point. This involves the computation $(D(x_1, x_2), D(x_1, x_3), \dots, D(x_1, x_n)), \dots, (D(x_i, x_{i+1}), D(x_i, x_{i+2}), \dots, D(x_i, x_n)), \dots, (D(x_{n-1}, x_n))$, which corresponds to an arithmetic series $n - 1 + n - 2 + \dots + 1$ of computations. Thus for agglomerative hierarchical clustering using *symmetrical* distances the number of distance computations is $n(n - 1)/2$ for the base level. At the next

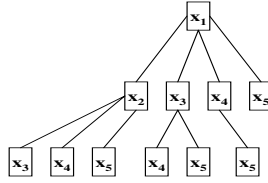


Fig. 4. A Simple Illustration for a Five Instance Problem of How the Triangular Inequality Can Save Distance Computations

level we need only recalculate the distance between the newly created cluster and the remaining $n - 2$ points and so on. Therefore, the total number of distance calculation is $n(n - 1)/2 + (n - 1)(n - 2)/2 = (n - 1)^2$. We can view the base level calculation pictorially as a tree construction as shown in Figure 4. If we perform the distance calculation at the first level of the tree then we can obtain bounds using the triangle inequality for **all** branches in the second level. This is as bounding the distance between two points requires the distance between these points and a common point, which in our case is x_1 . Thus in the best case there are only $n - 1$ distance computations instead of $(n - 1)^2$.

4.2 Average Case Analysis for Using the γ Constraint

However, it is highly unlikely that the best case situation will ever occur. We now focus on the average case analysis using the Markov inequality to determine the *expected* performance improvement which we later empirically verify. Let ρ be the average distance between any two instances in the data set to cluster. The triangle inequality provides a lower bound; if this bound exceeds γ , computational savings will result. We can bound how often this occurs if we can express γ in terms of ρ , hence let $\gamma = c\rho$.

Recall that the general form of the Markov inequality is: $P(X = x \geq a) \leq \frac{E(X)}{a}$, where x is a single value of the continuous random variable X , a is a constant and $E(X)$ is the expected value of X . In our situation since X is distance between two points chosen at random, $E(X) = \rho$ and $\gamma = a = c\rho$ as we wish to determine when the distance will exceed γ . Therefore, at the lowest level of the tree ($k = n$) then the *number of times* the triangle inequality will save us computation time is $n \frac{E(X)}{a} = n \frac{\rho}{c\rho} = n/c$, indicating a saving of a factor of $1/c$ at this lowest level. As the Markov inequality is a rather weak bound then in practice the saving may be substantially different as we shall see in our empirical section. The computation saving that are obtained at the bottom of the dendrogram are reflected at higher levels of the dendrogram. When growing the entire dendrogram we will save at least $n/c + (n - 1)/c \dots + 1/c$ distance calculations. This is an arithmetic sequence with the additive constant being $1/c$ and hence the total expected computations saved is at least $n/2(2/c + (n - 1)/c) = (n^2 + n)/2c$. As the total computations for regular hierarchical clustering is $(n - 1)^2$, the computational saving is expected to be by a approximately a factor of $1/2c$.

Consider the 150 instance IRIS data set ($n=150$) where the average distance (with attribute value ranges all being normalized to between 0 and 1) between two instances is 0.6; that is, $\rho = 0.6$. If we state that we do not wish to join clusters whose centroids are

Table 4. The Efficiency of Using the Geometric Reasoning Approach from Section 4 (Rounded Mean Number of Pair-wise Distance Calculations)

Data Set	Unconstrained	Using γ Constraint
Iris	22,201	19,830
Breast	487,204	431,321
Digit (3 vs 8)	3,996,001	3,432,021
Pima	588,289	501,323
Census	2,347,305,601	1,992,232,981
Sick	793,881	703,764

separated by a distance greater than 3.0, then $\gamma = 3.0 = 5\rho$. By not using the γ constraint and the triangle inequality the total number of computations is 22201, and the number of computations that are saved is at least $(150^2 + 150)/10 = 2265$; hence the saving is about 10%. We now show that the γ constraint can be used to improve efficiency of the basic agglomerative clustering algorithm. Table 4 illustrates the improvement that using a γ constraint equal to five times the average pairwise instance distance. We see that the average improvement is consistent with the average case bound derived above.

5 Constraints and Irreducible Clusterings

In the presence of constraints, the set partitions at each level of the dendrogram must be feasible. We have formally shown that if k_{max} is the maximum value of k for which a feasible clustering exists, then there is a way of joining clusters to reach another clustering with k_{min} clusters [5]. In this section we ask the following question: will traditional agglomerative clustering find a feasible clustering for each value of k between k_{max} and k_{min} ? We formally show that in the worse case, for certain types of constraints (and combinations of constraints), if mergers are performed in an arbitrary fashion (including the traditional hierarchical clustering algorithm, see Figure 1), then the dendrogram may prematurely dead-end. A premature dead-end implies that the dendrogram reaches a stage where no pair of clusters can be merged without violating one or more constraints, even though other sequences of mergers may reach significantly higher levels of the dendrogram. We use the following definition to capture the informal notion of a “premature end” in the construction of a dendrogram. How to perform agglomerative clustering in these dead-end situations remains an important open question.

Definition 3. A feasible clustering $C = \{C_1, C_2, \dots, C_k\}$ of a set S is **irreducible** if no pair of clusters in C can be merged to obtain a feasible clustering with $k - 1$ clusters.

The remainder of this section examines the question of which combinations of constraints can lead to premature stoppage of the dendrogram. We first consider each of the ML, CL and ϵ -constraints separately. It is easy to see that when only ML-constraints are used, the dendrogram can reach all the way up to a single cluster, no matter how mergers are done. The following illustrative example shows that with CL-constraints, if mergers are not done correctly, the dendrogram may stop prematurely.

Example: Consider a set S with $4k$ nodes. To describe the CL constraints, we will think of S as the union of four pairwise disjoint sets X, Y, Z and W , each with k nodes. Let $X = \{x_1, x_2, \dots, x_k\}$, $Y = \{y_1, y_2, \dots, y_k\}$, $Z = \{z_1, z_2, \dots, z_k\}$ and $W = \{w_1, w_2, \dots, w_k\}$. The CL-constraints are as follows. (a) There is a CL-constraint for each pair of nodes $\{x_i, x_j\}$, $i \neq j$, (b) There is a CL-constraint for each pair of nodes $\{w_i, w_j\}$, $i \neq j$, (c) There is a CL-constraint for each pair of nodes $\{y_i, z_j\}$, $1 \leq i, j \leq k$.

Assume that the distance between each pair of nodes in S is 1. Thus, nearest-neighbor mergers may lead to the following feasible clustering with $2k$ clusters: $\{x_1, y_1\}$, $\{x_2, y_2\}$, \dots , $\{x_k, y_k\}$, $\{z_1, w_1\}$, $\{z_2, w_2\}$, \dots , $\{z_k, w_k\}$. This collection of clusters can be seen to be irreducible in view of the given CL constraints. However, a feasible clustering with k clusters is possible: $\{x_1, w_1, y_1, y_2, \dots, y_k\}$, $\{x_2, w_2, z_1, z_2, \dots, z_k\}$, $\{x_3, w_3\}$, \dots , $\{x_k, w_k\}$. Thus, in this example, a carefully constructed dendrogram allows k additional levels. \square

When only the ϵ -constraint is considered, the following lemma points out that there is only one irreducible configuration; thus, no premature stoppages are possible. In proving this lemma, we will assume that the distance function is symmetric.

Lemma 1. *If S is a set of nodes to be clustered under an ϵ -constraint. Any irreducible and feasible collection C of clusters for S must satisfy the following two conditions.*

- (a) C contains at most one cluster with two or more nodes of S .
- (b) Each singleton cluster in C contains a node x with no ϵ -neighbors in S .

Proof: Suppose C has two or more clusters, say C_1 and C_2 , such that each of C_1 and C_2 has two or more nodes. We claim that C_1 and C_2 can be merged without violating the ϵ -constraint. This is because each node in C_1 (C_2) has an ϵ -neighbor in C_1 (C_2) since C is feasible and distances are symmetric. Thus, merging C_1 and C_2 cannot violate the ϵ -constraint. This contradicts the assumption that C is irreducible and the result of Part (a) follows. The proof for Part (b) is similar. Suppose C has a singleton cluster $C_1 = \{x\}$ and the node x has an ϵ -neighbor in some cluster C_2 . Again, C_1 and C_2 can be merged without violating the ϵ -constraint. \square

Lemma 1 can be seen to hold even for the combination of ML and ϵ constraints since ML constraints cannot be violated by merging clusters. Thus, no matter how clusters are merged at the intermediate levels, the highest level of the dendrogram will always correspond to the configuration described in the above lemma when ML and ϵ constraints are used. In the presence of CL-constraints, it was pointed out through an example that the dendrogram may stop prematurely if mergers are not carried out carefully. It is easy to extend the example to show that this behavior occurs even when CL-constraints are combined with ML-constraints or an ϵ -constraint.

6 Conclusion and Future Work

Our paper made two significant theoretical results. Firstly, the feasibility problem for *unspecified* k is studied and we find that clustering under all four types (ML, CL, ϵ and δ) of constraints is **NP**-complete; hence, creating a feasible dendrogram is intractable. These results are fundamentally different from our earlier work [4] because the feasibility problem and proofs are quite different. Secondly, we proved under some constraint

types (i.e. cannot-link) that traditional agglomerative clustering algorithms give rise to dead-end (irreducible) solutions. If there exists a feasible solution with k_{max} clusters then the traditional agglomerative clustering algorithm may not get all the way to a feasible solution with k_{min} clusters even though there exists feasible clusterings for each value between k_{max} and k_{min} . Therefore, the approach of joining the two “nearest” clusters may yield an incomplete dendrogram. How to perform clustering when dead-end feasible solutions exist remains an important open problem we intend to study.

Our experimental results indicate that small amounts of labeled data can improve the dendrogram quality with respect to cluster purity and “tightness” (as measured by the distortion). We find that the cluster-level δ constraint can reduce computational time between two and four fold by effectively creating a pruned dendrogram. To further improve the efficiency of agglomerative clustering we introduced the γ constraint, that allows the use of the triangle inequality to save computation time. We derived best case and expected case analysis for this situation which our experiments verified. Additional future work we will explore include constraints to create balanced dendrograms and the important asymmetric distance situation.

References

1. S. Basu, A. Banerjee, R. Mooney, Semi-supervised Clustering by Seeding, 19th *ICML*, 2002.
2. S. Basu, M. Bilenko and R. J. Mooney, Active Semi-Supervision for Pairwise Constrained Clustering, 4th *SIAM Data Mining Conf.*, 2004.
3. P. Bradley, U. Fayyad, and C. Reina, ”Scaling Clustering Algorithms to Large Databases”, 4th *ACM KDD Conference*. 1998.
4. I. Davidson and S. S. Ravi, ”Clustering with Constraints: Feasibility Issues and the k -Means Algorithm”, *SIAM International Conference on Data Mining*, 2005.
5. I. Davidson and S. S. Ravi, ”Towards Efficient and Improved Hierarchical Clustering with Instance and Cluster-Level Constraints”, Tech. Report, CS Department, SUNY - Albany, 2005. Available from: www.cs.albany.edu/~davidson
6. C. Elkan, Using the triangle inequality to accelerate k -means, *ICML*, 2003.
7. M. Garey and D. Johnson, *Computers and Intractability: A Guide to the Theory of NP-completeness*, Freeman and Co., 1979.
8. M. Garey, D. Johnson and H. Witsenhausen, ”The complexity of the generalized Lloyd-Max problem”, *IEEE Trans. Information Theory*, Vol. 28,2, 1982.
9. D. Klein, S. D. Kamvar and C. D. Manning, ”From Instance-Level Constraints to Space-Level Constraints: Making the Most of Prior Knowledge in Data Clustering”, *ICML* 2002.
10. M. Nanni, Speeding-up hierarchical agglomerative clustering in presence of expensive metrics, *PAKDD 2005*, LNAI 3518.
11. T. J. Schafer, ”The Complexity of Satisfiability Problems”, *STOC*, 1978.
12. K. Wagstaff and C. Cardie, ”Clustering with Instance-Level Constraints”, *ICML*, 2000.
13. D. B. West, *Introduction to Graph Theory*, Second Edition, Prentice-Hall, 2001.
14. K. Yang, R. Yang, M. Kafatos, ”A Feasible Method to Find Areas with Constraints Using Hierarchical Depth-First Clustering”, *Scientific and Stat. Database Management Conf.*, 2001.
15. O. R. Zaiane, A. Foss, C. Lee, W. Wang, On Data Clustering Analysis: Scalability, Constraints and Validation, *PAKDD*, 2000.
16. Y. Zho & G. Karypis, ”Hierarchical Clustering Algorithms for Document Datasets”, *Data Mining and Knowledge Discovery*, Vol. 10 No. 2, March 2005, pp. 141–168.

Cluster Aggregate Inequality and Multi-level Hierarchical Clustering

Chris Ding and Xiaofeng He

Lawrence Berkeley National Laboratory,
Berkeley, California 94720, USA

Abstract. We show that (1) in hierarchical clustering, many linkage functions satisfy a cluster aggregate inequality, which allows an exact $O(N^2)$ multi-level (using mutual nearest neighbor) implementation of the standard $O(N^3)$ agglomerative hierarchical clustering algorithm. (2) a desirable close friends cohesion of clusters can be translated into kNN consistency which is guaranteed by the multi-level algorithm; (3) For similarity-based linkage functions, the multi-level algorithm is naturally implemented as graph contraction. The effectiveness of our algorithms is demonstrated on a number of real life applications.

1 Introduction

Agglomerative hierarchical clustering (AHC) is developed in 1960's and is widely used in practice. AHC produces a tree describing the hierarchical cluster structure. Such a comprehensive description of the data is quite useful for broad areas of applications. For example, in bioinformatics research, AHC is most commonly used for clustering genes in a DNA gene microarray expression data, because the resulting hierarchical cluster structure is readily recognizable by biologists. The phylogenetic tree (similar to binary clustering tree) of organisms is often built using the UPGMA (unweighted pair group method average) AHC algorithm. In social sciences, the hierarchical cluster structure often reveals gradual evolving social relationships that help explain complex social issues. Another application of AHC is in classification tasks on a large dataset using support vector machine [12]. The hierarchical cluster structure allows one to use most detailed local representation near the decision boundaries where support vectors lie; but as one moves away from the decision boundaries, the centroid representation of progressively larger clusters can be used.

Besides hierarchical clustering, many other clustering algorithms have been developed (see recent survey and text books [5,1,3]). K -means clustering is perhaps the most commonly used method and is well developed. The gaussian mixture model using EM algorithm directly improves over the K -means method by using a probabilistic model of cluster membership of each object. Both algorithms can be viewed as a global objective function optimization problem. A related set of graph clustering algorithms are developed that partition nodes into two sub-clusters based on well-motivated clustering objective functions [8]. They

are typically applied in a top-down fashion (see also[7]), and thus complement the bottom-up AHC.

Standard AHC scales as $O(N^3)$. A detailed description of AHC and complete references can be found in [4,9]. A number of efficient implementation based on approximations have been proposed [10,6]. Several recent studies propose to integrate hierarchical clustering with additional information [2,11] or summary statistics[13].

In this paper, we focus on making AHC scale to large data set. Over the last 40 years, the basic AHC algorithm remains unchanged. The the basic algorithm is an iterative procedure; at each iteration, among all pairs of current clusters, the pair with largest linkage function value (smallest distance) are selected and merged.

We start with a key observation. In the AHC algorithm, at each iteration, one may merge all mutual-nearest-neighbor (1mn) pairs (defined by the linkage function) simultaneously in the same iteration. As long as the linkage function satisfies a “cluster aggregate inequality”, this modified algorithm of simultaneously merge all 1mn-pairs at each iteration produces identical clustering results as the standard AHC algorithm. The cluster aggregate inequality is satisfied by most common linkage functions (see §2.2). This modified algorithm provides a natural multi-level implementation of the AHC, where at each level we merge all 1mn pairs. This *multi-level hierarchical clustering* (MLHC) algorithm provides an order- N speedup over the standard AHC.

Next we propose “close friends” cohesion as a desirable feature for clustering, which requires that for every member in a cluster, its closest friend is also in the same cluster. We show that the MLHC guarantees the close-friends cohesion, a desirable feature for clustering. We further extend this cluster membership cohesion idea to (mutual) nearest neighbor consistency, and show that MLHC improves this KNN consistency compare to other clustering algorithms. (§3)

When the linkage function is expressed in similarity (in contrast to distance or dissimilarity), the new algorithm is identical to multi-level *graph contraction*. Graph contraction provides a natural framework for hierarchical clustering (§4).

The effectiveness of our algorithms is demonstrated on a number of real life applications: DNA gene expressions for lung cancer, global climate pattern, and internet newsgroups (§5).

2 Multi-level Hierarchical Clustering (MLHC)

2.1 Algorithm

The standard agglomerative hierarchical clustering is a bottom-up process. During each step, we merge two *mutual nearest neighbor* clusters C_p and C_q which are closest, or *mutual nearest neighbor* among all pairs of current clusters:

$$\min_{\langle pq \rangle} d(C_p, C_q).$$

where $d(\cdot, \cdot)$ is the dissimilarity-based (distance) linkage function between C_p and C_q . Many researches have studied the effects of different choice of linkage functions [4].

At each step of AHC with p current clusters, $p - 1$ new linkage functions need be computed, and we have total $O(p^2)$ pairwise linkage functions. It takes p^2 comparisons to search for the pair with the largest linkage. This is repeated $N - 1$ times. The total computation is

$$N_{\text{search}}^{\text{AHC}} = N^2 + (N - 1)^2 + \dots + 2^2 + 1^2 = O(N^3/3).$$

The new MLHC algorithm is motivated by the following observation. In each iterative step in AHC, when all pairwise linkage functions are computed, we can form all mutual nearest neighbor pairs (1mn-pairs) of current clusters using the linkage as the distance metric. Two objects (i, j) are a 1mn-pair if j is the nearest neighbor of i and vice versa.

In this perspective, the standard AHC merges only the 1mn-pair with largest linkage value. It is then natural to ask if we may also merge all other 1mn-pairs simultaneously. Will this produces the same results? An analysis shows that if the linkage function satisfies a ‘‘cluster aggregate inequality’’, then the clustering results remain the same as the standard AHC.

This observation suggests a simultaneous merging algorithm which we call MLHC. At each level, all 1mn-pairs are identified and merged (not just the pair with largest linkage value). This is repeated until all objects are merged into one cluster. The total number of level is about $\log_2 N$. Thus the required computation is approximately

$$N_{\text{search}}^{\text{MLHC}} = N^2 + (N/2)^2 + (N/4)^2 + 2^2 + 1^2 = O(4N^2/3).$$

The new algorithm speedups by a factor of order- N .

2.2 Cluster Aggregate Inequality

In this section, we show that the simultaneous merging of all 1mn-pairs in MLHC produces identical clustering results as the standard AHC, provided the linkage function satisfies a cluster aggregate inequality.

Definition. A linkage function $d(\cdot, \cdot)$ is said to satisfy the cluster aggregate inequality if for any three current clusters A, B, C . We try to merge B, C into a new cluster $B + C$. The cluster aggregate inequality is a property of the linkage function that the merged cluster $B + C$ is ‘‘no closer’’ to A than either one of its individual members B or C . More precisely, for distance (dissimilarity) based linkage $d(\cdot, \cdot)$, the cluster aggregate inequality is

$$d_{A, B+C} \geq \min(d_{A, C}, d_{A, B}) \tag{1}$$

for any triple (A, B, C) .

What kind of linkage functions satisfy the cluster aggregate inequality? It is interesting to see that most commonly used linkage functions satisfy cluster

aggregate inequality. Consider the four similarity-based linkage function. (i) the single linkage, defined as the closest distance among points in A, B , (ii) the complete linkage, defined as the farthest distance among points in A, B , (iii) the average linkage, defined as the average of all distances among points in A, B , (iv) the minimum variance linkage.

$$d_{A,B}^{\text{sgl}} = \min_{i \in A, j \in B} d_{ij} \quad (2)$$

$$d_{A,B}^{\text{cmp}} = \max_{i \in A, j \in B} d_{ij} \quad (3)$$

$$d_{A,B}^{\text{avg}} = \frac{d(A, B)}{n_A n_B} \quad (4)$$

$$d_{A,B}^{\text{min-var}} = \frac{n_A n_B}{n_A + n_B} \|\mathbf{c}_A - \mathbf{c}_B\|^2. \quad (5)$$

Theorem 1. The single link, the complete link and average link satisfy the strong cluster aggregate inequality.

Proof. For single link, one can easily see that

$$d_{A,B+C}^{\text{sgl}} = \min_{i \in A; j \in B+C} d_{ij} = \min(\min_{i \in A; j \in B} d_{ij}, \min_{i \in A; j \in C} d_{ij}) = \min(d_{A,B}^{\text{sgl}}, d_{A,C}^{\text{sgl}})$$

Thus the equality in Eq.(1) hold for single link. With same reasoning, one can see the inequality holds for complete linkage.

For average link, we have

$$\begin{aligned} d_{A,B+C}^{\text{avg}} &= \sum_{i \in A; j \in B+C} \frac{d_{ij}}{|A||B+C|} \\ &= \sum_{i \in A; j \in B} \frac{d_{ij}}{|A||B+C|} + \sum_{i \in A; j \in C} \frac{d_{ij}}{|A||B+C|} \\ &= \frac{|B|}{|B+C|} d_{A,B}^{\text{avg}} + \frac{|C|}{|B+C|} d_{A,C}^{\text{avg}} \\ &\geq \frac{|B|}{|B+C|} \min(d_{A,B}^{\text{avg}}, d_{A,C}^{\text{avg}}) + \frac{|C|}{|B+C|} \min(d_{A,B}^{\text{avg}}, d_{A,C}^{\text{avg}}) \\ &= \min(d_{A,B}^{\text{avg}}, d_{A,C}^{\text{avg}}) \quad \square \end{aligned}$$

2.3 Equivalence of MLHC and AHC

Cluster aggregate inequality plays a fundamental role in hierarchical clustering. It is similar to the triangle inequality in metric space: for any three vectors $\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k$ in the Hilbert space with the distance metric $d(\cdot, \cdot)$, the metric must satisfies the triangle inequality

$$d_{ik} \leq d_{ij} + d_{jk}.$$

Triangle inequality plays a fundamental role in deriving properties of the metric Space. Now we prove the main results of this paper.

Theorem 2. If the linkage function satisfies the cluster aggregate inequality, the clustering trees produced by MLHC is identical to that produced by standard AHC.

Proof. We view the linkage function as a distance metric and build all 1mn-pairs based on the linkage function. We show that the AHC is iteratively merging 1mn-pairs, which is same as MLHC. The details is broken into two features of AHC below. \square

We first note a simple feature of AHC:

Feature 1. The closest pair must be a 1mn-pair.

Proof. Suppose this is not true, i.e., the closest pair is (a, b) , but a is not the closest neighbor of b . There must exist another point c which is the closest neighbor of a . Then the pair (a, c) must be the closest pair, but this contradicts the fact that (a, b) is the closest pair. Thus a must be the closest neighbor of b . Similarly, b must be the closest neighbor of a . \square

Next, we prove a key feature of AHC. This shows the essence of Theorem 2.

Feature 2. If the linkage function satisfies the cluster aggregate inequality, then any 1mn-pair must be preserved and will merge eventually in standard AHC.

Proof. Suppose at certain iteration, the current clusters are listed as $(C_{j_1}, C_{j_2}), (C_{j_3}, C_{j_4}), C_{j_5}, (C_{j_6}, C_{j_7}), \dots$, where 1mn-pair is indicated by the parenthesis. In AHC, the 1mn-pair with largest linkage value, say (C_{j_6}, C_{j_7}) , is merged. Due to the cluster aggregate inequality, the newly merged cluster $C_{(j_6, j_7)}$ will be “no closer” to any other current clusters. Thus the 1mn of $C_{(j_6, j_7)}$ can not be any member of the current remaining 1mn-pairs, say C_{j_1} .

This can seen as follows. By construction, neither C_{j_6} nor C_{j_7} is closer to C_{j_1} than C_{j_2} does. Due to the cluster aggregate inequality, the newly merged cluster $C_{(j_6, j_7)}$ will be “no closer” to C_{j_1} than either C_{j_6} or C_{j_7} does. Thus 1mn of $C_{(j_6, j_7)}$ can not be C_{j_1} .

This guarantees that the 1mn-pair (C_{j_1}, C_{j_2}) will never be broken by a merged cluster. Thus in the next iteration, either a current 1mn-pair, say (C_{j_3}, C_{j_4}) , is merged, or the newly-merged $C_{(j_6, j_7)}$ is merged with a singleton cluster, say C_{j_5} . Therefore, the 1mn-pair (C_{j_1}, C_{j_2}) will preserve and eventually merge at some later iteration. \square

We give a simple example to illustrating some of the concepts. In Figure 1(b), we have 5 objects. They form two 1mn-pairs (a, b) , (d, e) and one isolated object c . We do the standard AHC. Suppose (a, b) has the largest linkage value. So a, b are first merged into $(a+b)$. We assert that the 1mn-pair (d, e) must be preserved and will merge in later stages in AHC. This is done in two stages. First, we show that d cannot be the nearest neighbor of $(a+b)$, i.e.,

$$d(d, a+b) > d(d, e). \quad (6)$$

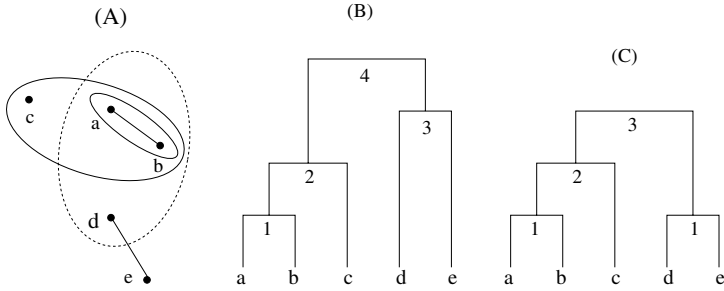


Fig. 1. (A) Dataset of 5 objects. 1mn-pairs are linked by a line. A merge is indicated by an elongated solid circle. (B) Dendrogram using AHC. Numbers indicate order of the merge. (C) Dendrogram using MLHC.

The fact that (d, e) is a 1mn-pair implies

$$d(d, a) > d(d, e), \tag{7}$$

$$d(d, b) > d(d, e), \tag{8}$$

$$d(d, c) > d(d, e). \tag{9}$$

From the cluster aggregate inequality, $d(d, a + b) \geq \min[d(d, a), d(d, b)]$. From this, and making use of Eqs.(7, 8), we obtain Eq.(6).

In the next round of AHC, either (i) the pair $(a + b, c)$ has the largest linkage value, or (ii) the pair (d, e) has the largest linkage value. If (ii) holds, our assertion is proved. Suppose (i) holds. Then $(a + b, c)$ are merged into $(a + b + c)$. Now we show that d cannot be the nearest neighbor of $(a + b + c)$, i.e.,

$$d(d, a + b + c) > d(d, e). \tag{10}$$

From the cluster aggregate inequality, $d(d, a + b + c) \geq \min[d(d, a + b), d(d, c)]$. From this, together with Eqs.(6, 9), we obtain Eq.(10). Therefore, (d, e) is the pair with the largest linkage value. Thus (d, e) are merged into $(d + e)$.

This example shows how 1mn pairs are preserved in AHC. Thus any cluster merge in MLHC (those 1mn-pairs) will also occur in AHC. There are total $N - 1$ cluster merges in both MLHC and AHC. So any cluster merge in AHC also occur in MLHC.

Both AHC and MLHC algorithms lead to the same binary cluster tree. The difference is the sequential order they are merged. This is illustrated in Figure 1.(B,C). If we represent the tree height by the standard linkage function value for each merge, the dendrograms of the two algorithm remains identical.

We emphasize that the equivalence of MLHC and AHC only requires 1mn-pair preservation during AHC, as shown in Feature 2 in the above. Therefore, cluster aggregate inequality is a sufficient condition for 1mn-pair preservation.

For a given dataset, it is possible that a particular linkage function maybe not satisfy the generic cluster aggregate inequality for all possible triples (i, j, k) ,

but the 1mn-pair preservation holds during AHC and thus MLHC is equivalent to AHC.

In summary, these analysis not only shows that MLHC is equivalent to AHC, thus providing a $O(n)$ speedup; but also brought out a new insight for AHC, i.e., 1mn-pair preservation during AHC. This leads to close-friends cohesion.

3 “Close Friends” Cohesion

One of the fundamental concept of data clustering is that members of the same cluster have high association with each other. One way to characterize the within cluster association is the cohesion of the cluster members via the preservation of “close friends”. Suppose we divide 100 people into 10 teams. The cohesion of each team is greatly enhanced if for any member in a team, his/her close friends are also in the same team.

It is interesting to note that this close-friends cohesion is guaranteed by MLHC, if we interpret 1mn-pair relationship as close friends; by construction, this cohesion is guaranteed at all levels. We say the clustering results satisfy 1mn-consistency if for each object in a cluster, its 1mn is also in the same cluster.

By Theorem 2, clusters produced by the standard AHC also enjoy the same 1mn-consistency as in MLHC. We summarize this important result as

Theorem 3. In MLHC, 1mn-consistency is fully guaranteed. In agglomerative hierarchical clustering, if the linkage function satisfy the cluster aggregate inequality, 1mn-consistency is fully guaranteed.

3.1 Cluster Membership kNN Consistency

1mn describe the “most close” friend. Beyond the closest friend, it is desirable that other less close friends are also in the same cluster. This will increase cohesiveness of the clustering.

We thus further extend the “close friends” into k-nearest-neighbor, the cohesiveness of a cluster becomes the following knn consistency:

Cluster Membership k NN Consistency: For any data object in a cluster, its k -nearest neighbors are also in the same cluster.

Note that the relationship of “close friend” is not symmetric. Although a 's closest friend is b , b 's closest friend could be c , not a . Thus the “mutual closest friend” implies the tightest friendship. Thus k -Mutual-Nearest-Neighbor Consistency is more desirable.

3.2 Enforcing kNN Consistency

In general, clustering algorithms perform global optimizations, such as K -means algorithm, will gives final clusters with a smaller degree of cluster knn and kmn consistency than the bottom hierarchical clustering. This is because that in global optimizations, nearest-neighbor relations are not considered explicitly.

In HC, a pair of clusters are merged if their linkage value is high. Thus clusters being merged are typically very similar to each other. Therefore the nearest-neighbor local information are utilized to some extent, leading to higher degree of knn consistency.

What about other knn/kmn consistency? First we note that cluster knn consistency defines a “transitive relationship”. If x_1 is a member of cluster C_1 , and x_2 is the 1nn of x_1 , then by cluster 1nn consistency, x_2 should also be a member of C_1 . This transitive relation implies that 100% 1nn consistency can be achieved only if entire connected component of the 1nn graph are in C_1 . To generate clusters that guarantee 100% knn consistency, at each level of the MLHC, we must first generate knn-graph, identify all connected components, and for each CC, merge all current clusters into one cluster.

Because for any object, its 2nn set always include its 1nn set. Thus 2nn consistency guarantees 1nn consistency. Similarly, because any knn set include kmn set, knn consistency implies (k-1)-nn consistency.

4 Similarity-Based Hierarchical Clustering: Multi-level Graph Contraction

The above discussion uses the distance-based linkage. All results there can be easily translate into similarity-based linkage function.

For similarity-based linkage we select the pair with the largest linkage to merge: $\max_{\langle pq \rangle} s(C_p, C_q)$, where $s(C_p, C_q)$ is the aggregate similarity between clusters C_p, C_q . Let the initial pairwise similarity are $W = (w_{ij})$. The aggregate similarity has a simple form, $s(C_p, C_q) = \sum_{i \in C} \sum_{j \in C} w_{ij}$.

Cluster aggregate inequality using similarity-based linkage can be written as

$$s(A, B + C) \leq \max[s(A, B), s(A, C)] \tag{11}$$

Consider the following similarity-based linkage functions. (i) the single linkage, defined as the largest similarity among points in A, B , (ii) the complete linkage, defined as the smallest similarity among points in A, B , (iii) the average linkage, defined as the average of all similarities among points in A, B ,

$$s_{\text{single}}(A, B) = \max_{i \in A} \max_{j \in B} s_{ij}, \quad s_{\text{complete}}(A, B) = \min_{i \in A} \min_{j \in B} s_{ij}, \quad s_{\text{avg}}(A, B) = \frac{s(A, B)}{|A||B|}.$$

With similar analysis as in the case of distance-based clustering, we can proof

Theorem 4. All three above similarity-based linkage functions satisfy the cluster aggregate inequality. The similarity based linkage functions have an advantage that merging two cluster become graph node contraction. Defining the similarity between two objects as the weight on an edge between them, this forms a similarity graph. Thus the multi-level hierarchical clustering naturally become multi-level graph contraction of the similarity graph. Many well-known results in graph theory can be applied.

Merging two current clusters into a new cluster corresponds to contracting two nodes i, j into a new node k and with edge e_{ij} being eliminated. Weights of the graph are updated according to standard graph contraction procedure. Let $W^{(t)}, W^{(t+1)}$ be the weights of the similarity graph at steps $t, t + 1$. The updated weights for contracting the edge e_{ij} and merging nodes i, j into k are

$$\begin{cases} w_{kk}^{(t+1)} = w_{ii}^{(t)} + w_{jj}^{(t)} + w_{ij}^{(t)} \\ w_{kp}^{(t+1)} = w_{ip}^{(t)} + w_{jp}^{(t)}, \quad \forall p \notin \{i, j, k\} \\ w_{pq}^{(t+1)} = w_{pq}^{(t)}, \quad \forall p, q \notin \{i, j, k\} \end{cases}$$

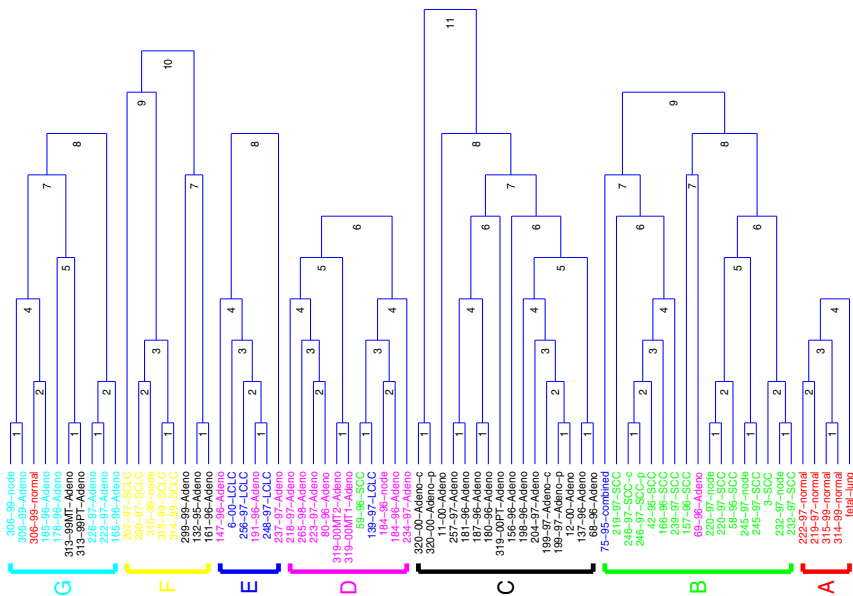


Fig. 2. MLHC clustering of lung cancer gene expressions

Using the above graph contraction operation, AHC and MLHC can be described very succinctly. AHC can be viewed as contracting one edge (with heaviest weight) at each step. The steps are repeated until all nodes merge into one. This takes exactly $N - 1$ steps. MLHC can be viewed as contracting all $1mn$ -pairs simultaneously at each step. It is repeated about $O(\log_2 N)$ times. Now the speedup of MLHC over AHC is clear. At each step, all-pair linkage function computation is necessary. But the number of steps required in AHC is $N - 1$, and it is $O(\log_2 N)$ in MLHC.

5 Experiments

5.1 DNA Gene Expression for Lung Cancer

The DNA gene expressions of lung cancer patients (available online: http://genome-www.stanford.edu/lung_cancer/adeno/) contains 73 samples of 67 lung tumors from patients whose clinical course was followed for up to 5 years. The samples comprise of 916 DNA clones representing 835 unique genes. The samples are classified into 5 groups by visual examination (41 Adenocarcinomas (ACs), 16 squamous cell carcinomas (SCCs), 5 large cell lung cancers(LCLCs), 5 small cell lung cancer (SCLCs) and 5 normal tissue with one fetal lung tissue.). The largest group ACs is further classified into three smaller groups. The purpose is to see if we can recover this 7 groups using unsupervised learning method, i.e., the hierarchical clustering method. The Pearson correlations c_{ij} among tissue samples are computed first, and the similarity metric is defined as $s_{ij} = \exp(5c_{ij})$. We use MLHC and obtain the cluster structure as shown in Fig.2.

As the Figure shows, at 1st level, after an all-pair computation, 18 1mn-pairs are formed and merged. At 2nd level, 11 1mn-pairs are formed and merged. Total 11 levels of merges are required to obtain 7 clusters. In contrast, for standard AHC, We need 66 levels of merge steps to obtain 7 clusters. The clustering

result is give in the confusion matrix $= \begin{bmatrix} 5 & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & 15 & \cdot & 1 & 1 & \cdot \\ \cdot & \cdot & 16 & \cdot & \cdot & \cdot \\ \cdot & 1 & \cdot & 9 & 1 & \cdot \\ \cdot & \cdot & \cdot & 3 & 3 & \cdot \\ \cdot & \cdot & \cdot & 3 & \cdot & 5 \\ 1 & \cdot & 2 & \cdot & \cdot & 7 \end{bmatrix}$ where $T = (t_{ij})$, t_{ij} is the

number of data points which are observed to be in cluster i , but was computed via the clustering method to belong to cluster j . The accuracy is defined as $Q = \sum_k t_{kk}/N = 82\%$. indicating the effectiveness of the clustering algorithm.

5.2 Climate Pattern

We tested MLHC on global precipitation data as shown in Fig.3. Regularly-spaced data points cover the surface of earth. Each data point is a 402-dimensional vector containing seasonal means over 100 years and geometric information: longitude and latitude. Similarity between two points are based two factors: (1) precipitation pattern similarity computed as Euclidean distance and (2) geometric closeness based on simple physical distance. The obtained stable regions (shown in different color and symbols) correlate well with continents, say, in Australia, south Americas.

5.3 Internet Newsgroups

We apply MLHC on Internet newsgroup articles. A 20-newsgroup dataset is from www.cs.cmu.edu/~afs/cs/project/theo-11/www/naive-bayes.html. 1000 words are selected according to the mutual information between words and documents in unsupervised manner. Word - document matrix is first constructed using standard `tf.idf` term weighting. Cosine similarity between documents is used. We focus on two sets of 5-newsgroup combinations listed below:

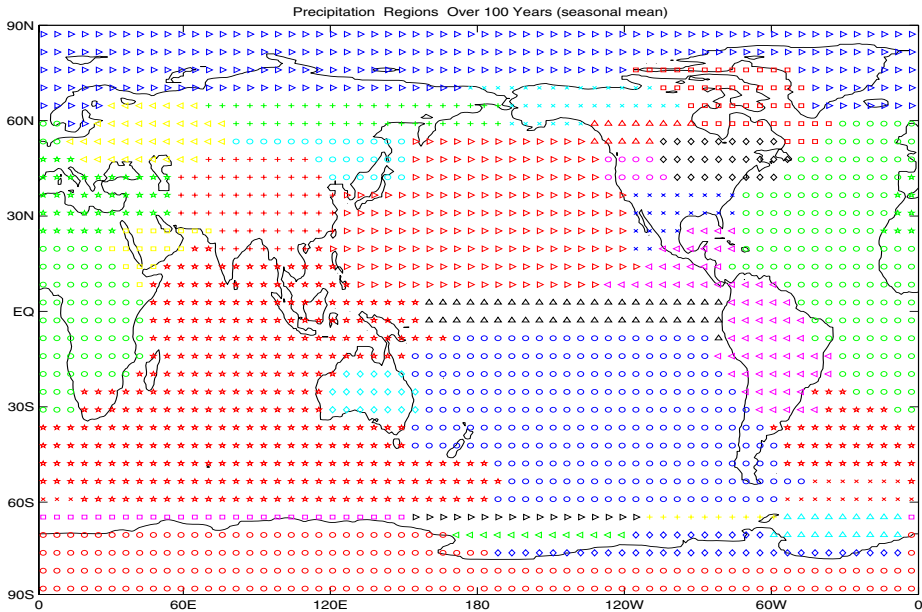


Fig. 3. Global precipitation pattern based on seasonal means over 100 years

A5:	B5:
NG2: comp.graphics	NG2: comp.graphics
NG9: rec.motorcycles	NG3: comp.os.ms-windows
NG10: rec.sport.baseball	NG8: rec.autos
NG15: sci.space	NG13: sci.electronics
NG18: talk.politics.mideast	NG19: talk.politics.misc

In A5, clusters overlap at medium level. In B5, clusters overlap substantially. Table 1 contains the results of MLHC. To accumulate sufficient statistics, for each newsgroup combination, we generate 10 datasets, each of which is a random sample of documents from the newsgroups (with 100 documents per newsgroup). The results in the table are the average over these 10 random sampled datasets. For comparison purpose, we also run K -means clustering. For each dataset, we run 10 K -means clustering from random seeds for cluster centroids and select the best result as determined by the K -means objective function value. Results are given in Table 1. Note that because percentage consistency results are close to 1, we give inconsistency = 1 - consistency.

From Table 1, the MLHC results have better clustering accuracy (last column of Table 1) compared to K -means clustering. More important is MLHC always provides clustering with better kmn cluster membership consistency. For 1nn-consistency, MLHC is perfect since this is guaranteed by MLHC. With this, it is not surprising that MLHC has substantially better 1nn-consistency than K -means method, about half as smaller. In all categories, MLHC has better kmn/kmn consistency than K -means.

Table 1. Fractional knn and kmn inconsistency and clustering accuracy (last column) for newsgroup datasets A5 and B5. For dataset A5, 1nn inconsistency is 16.2% for K-means and 8.5% for MLHC.

	1nn	2nn	3nn	1mn	2mn	3mn	Accuracy
A5							
K-means	16.2	28.4	37.8	6.4	14.5	23.0	75.2%
MLHC	8.5	24.1	36.4	0	6.9	16.6	77.6%
B5							
K-means	23.1	39.4	50.6	8.5	21.6	32.8	56.3%
MLHC	10.2	28.9	45.0	0	9.3	21.3	60.7%

6 Summary

In this paper, we propose a modification of the standard AHC algorithm that allow an order- N faster implementation. The modification is based on the recognition that all 1mn-pairs in each iteration of AHC can be merged if the linkage function satisfies the cluster aggregate inequality. This leads to the multi-level hierarchical clustering algorithm. Many commonly used linkage functions satisfy this inequality and thus will benefit from this modification. We propose “close friends” cohesion as important feature of clustering and show that it is fully guarantees in the algorithm. This is further extended to cluster membership KNN consistency. Experiments on newsgroup show that kNN consistency is satisfied much better by MLHC than widely used algorithms such as K -means .

References

1. R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification, 2nd ed.* Wiley, 2000.
2. B. Fung, K. Wang, and M. Ester. Large hierarchical document clustering using frequent itemsets. *Proc. SIAM Data Mining Conf*, 2003.
3. T. Hastie, R. Tibshirani, and J. Friedman. *Elements of Statistical Learning*. Springer Verlag, 2001.
4. A.K. Jain and R.C. Dubes. *Algorithms for clustering data*. Prentice Hall, 1988.
5. A.K. Jain, M.N. Murty, and P.J. Flynn. Data clustering: a review. *ACM Computing Surveys*, 31:264–323, 1999.
6. S.Y. Jung and T.-S. Kim. An agglomerative hierarchical clustering using partial maximum array and incremental similarity computation method. *Proc. SIAM Conf. on Data Mining*, pages 265–272, 2001.
7. G. Karypis, E.-H. Han, and V. Kumar. Chameleon: Hierarchical clustering using dynamic modeling. *IEEE Computer*, 32:68–75, 1999.
8. J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE. Trans. on Pattern Analysis and Machine Intelligence*, 22:888–905, 2000.
9. S. Theodoridis and K. Koutroubas. *Pattern Recognition*. Academic Press, 1999.

10. E.M. Voorhees. Implementing agglomerative hierarchic clustering algorithms for use in document retrieval. *Information Processing and Management*, 22:465–476, 1986.
11. H Xiong, M. Steinbach, P-N. Tan, and V. Kumar. Hicap:hierarchical clustering with pattern preservation. *Proc. SIAM Data Mining Conf*, pages 279–290, 2004.
12. H. Yu, J. Yang, and J. Han. Classifying large data sets using svms with hierarchical clusters. In *Proc. ACM Int'l Conf Knowledge Disc. Data Mining (KDD)*, pages 306–315, 2003.
13. T. Zhang, R. Ramakrishnan, and M. Livny. Birch: an efficient data clustering method for very large databases. *Proc. ACM Int'l Conf. Management of Data (SIGMOD)*, pages 103–114, 1996.

Ensembles of Balanced Nested Dichotomies for Multi-class Problems

Lin Dong¹, Eibe F. Franke¹, and Stefan Kramer²

¹ Department of Computer Science, University of Waikato, New Zealand
{ld21, eibe}@cs.waikato.ac.nz

² Department of Computer Science, Technical University of Munich, Germany
kramer@in.tum.de

Abstract. A system of nested dichotomies is a hierarchical decomposition of a multi-class problem with c classes into $c - 1$ two-class problems and can be represented as a tree structure. Ensembles of randomly-generated nested dichotomies have proven to be an effective approach to multi-class learning problems [1]. However, sampling trees by giving each tree equal probability means that the depth of a tree is limited only by the number of classes, and very unbalanced trees can negatively affect runtime. In this paper we investigate two approaches to building balanced nested dichotomies—*class-balanced* nested dichotomies and *data-balanced* nested dichotomies—and evaluate them in the same ensemble setting. Using C4.5 decision trees as the base models, we show that both approaches can reduce runtime with little or no effect on accuracy, especially on problems with many classes. We also investigate the effect of caching models when building ensembles of nested dichotomies.

1 Introduction

Many real-world classification problems are multi-class problems: they involve a nominal class variable having more than two values. The traditionally used approach for tackling this type of problem is One-vs-All the learning algorithm to deal with multi-class problems directly, and the other is one-vs-one level-wise classification and forming multi-classification based on the decision obtained from the two-class problem. The latter approach is appealing because it does not involve any change of the underlying two-class learning algorithm. Well-known examples of this type of approach are *one-vs-one* coding [2] and *pairwise classification* [3], and they often result in significant increases in accuracy.

Recently, it has been shown that an ensemble of nested dichotomies is a promising alternative to pairwise classification and *one-vs-one* coding. In experiments with a decision tree learner and logistic regression, the performance was significantly better than the base learner used, and they yield a probability estimate in a natural and well-founded way if the base learner can generate two-class probability estimates [1].

Adapting an ensemble of nested dichotomies, a learner used for pairwise classification, to the multi-class problem can increase its performance. Although pairwise classification is a promising approach to the base learner, $c * (c - 1) / 2$ is the number of

with each class, each learning problem is challenging, although the original problem became only data for the relevant pair of classes considered [3]. Applying a learning algorithm has scale linearly in the number of instances, and applying has every class has the same number of instances, the overall number of pairs with each class is linear in the number of classes.¹

Building a single ensemble of nested dichotomies in the above learning algorithm is linear in the number of classes in the worst case, but the algorithm can be applied a fixed, pre-specified number of times to build an ensemble of size (10 to 20 ensemble members), we found to be generally sufficient to achieve maximum accuracy on the UCI data sets investigated in [1].

In this paper, we are looking at a approach of combining the ensemble to build an ensemble of nested dichotomies (END). More specifically, we consider two different approaches of nested dichotomies (ECBND and EDBND, respectively). Using C4.5 as the base learner, we show how they can improve performance, especially on problems with many classes, with little or no effect on classification accuracy. We also investigate the effect of caching models: the above two-class learning problem can occur multiple times in an ensemble and it is an ensemble of two-class base models that have already been built for previous ensemble of nested dichotomies.

The paper is organized as follows. In Section 2 we discuss the basic methodology of building END, our modified version of the algorithm (ECBND and EDBND), and the use of caching models. Section 3 presents the empirical results, obtained from 21 multi-class UCI datasets, and several artificial domains with a varying number of classes. Section 4 concludes and points to future work.

2 Balanced Nested Dichotomies

Any ensemble of nested dichotomies is a practical model has a reduced number of attributes, one class in one or more classes (e.g. [4] in decision tree). The decomposition can be represented as a binary tree (Figure 1). Each node of the tree is a pair of class labels, the corresponding training data and a binary classifier. At the very beginning, the root node contains the whole set of the original class labels, corresponding to the multi-class classification problem. This is then split into two subsets. The two subsets of class labels are each a two "leaf" class and a binary classifier, learned from training the set. The training data is split into two subsets corresponding to the two leaf classes and one subset of training data is regarded as the positive example while the other subset is regarded as the negative example. The root node of the tree in the worst case of the original class label with the corresponding training data and a leaf node is built by applying this procedure recursively. The procedure finally reaches a leaf node if the node contains only one class label.

¹ Pairwise classification is actually even more beneficial when the base learner's runtime is worse than linear in the number of instances.

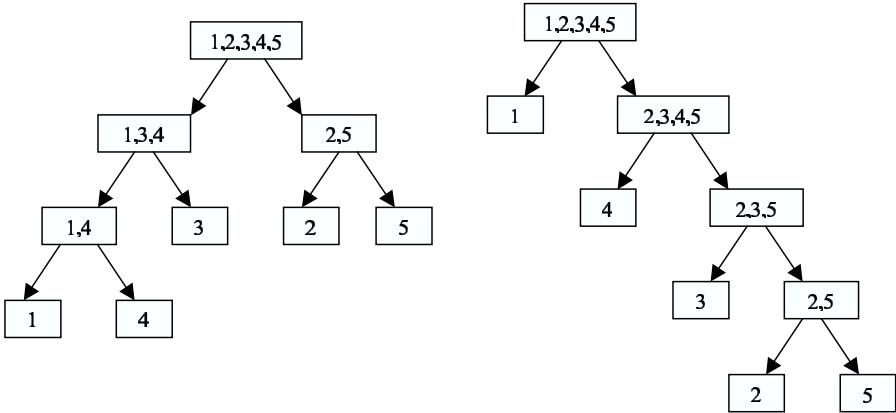


Fig. 1. Two different nested dichotomies for a classification problem with five classes

It is obvious that for any given c -class problem, there exist c leaf nodes (one for each class) and $c-1$ internal nodes. Each internal node contains a binary classification. A nice feature of using a sequence of nested dichotomies for a c -class problem is that they yield classification probabilities that are in a straightforward fashion. Assuming the individual wo-class classification is handled by the base learner, once a classification probability is available for the corresponding wo-class problem—one for each branch extending from the corresponding internal node—we can obtain a classification probability for a specific leaf node (i.e. the class label in the c -class problem) has associated with it leaf node) by simply multiplying together the probabilities that are obtained for the binary classification along the specific path from the root node to the leaf [1].

However, there is a problem with the application of nested dichotomies to standard c -class problems: there are any number of ways to choose a given set of classes, and in the absence of prior knowledge about which the specific decomposition is more appropriate, it is not clear which one to use. Figure 1 shows two different ways of nested dichotomies for a five-class problem. The problem is that two different ways will often lead to different predictions because the binary classification for each node are dealing with different wo-class problems. The selection of the sequence of classes will influence the classification result. Given this observation and the success of ensemble learning in yielding accurate predictions, it is a reasonable question to ask whether there are any other possible nested dichotomies for a given problem and average their classification probabilities. Unfortunately this is infeasible because the number of possible ways of nested dichotomies for a c -class problem is $(2c-3)!!$ [1]. Hence it is necessary to have a better way.

Franz and Kramer [1] also found only for the case of all possible ways, giving each way equal probability. However, it is not clear whether this will be a good approach. In the absence of prior knowledge, any algorithm that has no given preference for a particular class can be considered as a suitable candidate. The problem with the algorithm based on a uniform distribution

Table 1. Comparison of the number of possible trees

Number of classes	Number of nested dichotomies	Number of class-balanced nested dichotomies
2	1	1
3	3	3
4	15	3
5	105	30
6	945	90
7	10,395	315
8	135,135	315
9	2,027,025	11,340
10	34,459,425	113,400
11	654,729,075	1,247,400
12	13,749,310,575	3,742,200

over, see it has the tree depth is only limited by the number of classes, and deeper trees can achieve a long linear build. Consider the case where the tree is a linear tree in the second tree shown in Figure 1, and assume the worst case class are equally sized (class 2 and 5 in the example). Then all binary learning problems will involve the worst case class, increasing a high combinatorial cost for the process of building the binary classifier in the tree.

2.1 Class-Balanced Nested Dichotomies

In light of the evaluation we consider two different learning strategies in this paper. The first method is based on balancing the number of classes at each node. Instead of learning for the space of all possible trees, we search for the space of all balanced trees, and build an ensemble of balanced trees. The advantage of this method is that the depth of the tree is guaranteed to be logarithmic in the number of classes. We call this an ensemble of class-balanced nested dichotomies (ECBND).

The number of possible class-balanced nested dichotomies is obviously smaller than the total number of nested dichotomies. The following recurrence relation defines the number of possible class-balanced trees:

$$T(c) = \begin{cases} \frac{1}{2} \binom{c}{c/2} T(\frac{c}{2})T(\frac{c}{2}) & : \text{ if } c \text{ even} \\ \binom{c}{(c+1)/2} T(\frac{c+1}{2})T(\frac{c-1}{2}) & : \text{ if } c \text{ odd} \end{cases}$$

where $T(1) = 1$ and $T(2) = 1$.

Table 1 shows the number of possible binary trees of nested dichotomies for up to 12 classes for the class-balanced (CBND) and the unconstrained case (ND). It shows that a non-trivial number of CBND can be generated for classification problems with very few classes.

Figure 2 shows the algorithm for building a binary tree of class-balanced nested dichotomies. At each node the set of classes is linearly ordered by size (of course, if the number of classes is odd, there will not be exactly equal), and the

```

method buildClassBalancedNestedDichotomies(Dataset  $D$ , Set of classes  $C$ )
    if  $|C| = 1$  then return
     $P$  = subset of  $C$ , randomly chosen from all subsets of size  $\lfloor |C|/2 \rfloor$ 
     $N = C \setminus P$ 
     $D_p$  = all instances in  $D$  apart from those pertaining to classes in  $P$ 
    buildClassBalancedNestedDichotomies( $D_p$ ,  $P$ )
     $D_n$  = all instances in  $D$  apart from those pertaining to classes in  $N$ 
    buildClassBalancedNestedDichotomies( $D_n$ ,  $N$ )
     $D'$  = a two-class version of  $D$  created based on  $N$  and  $P$ 
    classifierForNode = buildClassifier( $D'$ )
    
```

Fig. 2. Algorithm for generating class-balanced nested dichotomies

balance learning algorithm is applied on the data according to the workflow. The algorithm then recursively splits only one class left. It is applied repeatedly with different and non-sequential nodes to generate a collection of nodes.

2.2 Data-Balanced Nested Dichotomies

The essential problem with the class-balanced approach: one class imbalance is always unbalanced and one class is always chosen to be the majority. In this case a class-balanced node does not only have a majority class (i.e. has an imbalance of instances in the majority node of an internal node). This can negatively affect, namely if the balance learning algorithm has a complexity worse than linear in the number of instances. Hence we also consider a simple algorithm for building data-balanced nested dichotomies in this case. Note that this method violates the condition that the algorithm should not be biased towards a specific class: balanced on the child node for large class will be located higher in the tree structure. Despite this problem we decided to investigate this method because it is difficult to say how important this condition is in practice.

Figure 3 shows how our algorithm for building a type of data-balanced nested dichotomies. It randomly assigns classes to workflow until the size of the remaining data in one of the branches exceeds half the total amount of remaining data at the node. One motivation for using this simple algorithm was that it is important to maintain a degree of randomness in the assignment of classes to branches in order to avoid diversity in the collection of randomly generated types of nested dichotomies. Given that we are aiming for an ensemble of nested dichotomies it would not be advisable, even if it were computationally feasible, to aim for an optimal balance because this would eventually result in the number of nodes that can be generated. Even with our simple algorithm diversity is affected when the class distribution is very unbalanced. However, it is difficult to derive a general exception for the number of nodes that can be generated by this method because this number depends on the class distribution in the data set.

```

method buildDataBalancedNestedDichotomies(Dataset  $D$ , List of classes  $C$ )

  if  $|C| = 1$  then return
   $C =$  random permutation of  $C$ 
   $D_p = \emptyset, D_n = \emptyset$ 
  do
    if ( $|C| > 1$ ) then
      add all instances from  $D$  pertaining to first class in  $C$  to  $D_p$ 
      add all instances from  $D$  pertaining to last class in  $C$  to  $D_n$ 
      remove first and last class from  $C$ 
    else
      add all instances from  $D$  pertaining to remaining class in  $C$  to  $D_p$ 
      remove remaining class from  $C$ 
  while ( $|D_p| < \lfloor |D|/2 \rfloor$ ) and ( $|D_n| < \lfloor |D|/2 \rfloor$ )
  if ( $|D_p| \geq \lfloor |D|/2 \rfloor$ ) then
    add instances from  $D$  pertaining to remaining classes in  $C$  to  $D_n$ 
  else
    add instances from  $D$  pertaining to remaining classes in  $C$  to  $D_p$ 
   $P =$  all classes present in  $D_p, N =$  all classes present in  $D_n$ 
  buildDataBalancedNestedDichotomies( $D_p, P$ )
  buildDataBalancedNestedDichotomies( $D_n, N$ )
   $D' =$  a two-class version of  $D$  created based on  $N$  and  $P$ 
  classifierForNode = classifier learned by base learner from  $D'$ 

```

Fig. 3. Algorithm for generating data-balanced nested dichotomies

2.3 Computational Complexity

The motivation for using balanced nested dichotomies is highlighted in Section 2.1. In the following we analyze the computational complexity of completely and balanced nested dichotomies. Let c be the number of classes in the data set, and n be the number of training instances. For simplicity, assume all classes have an approximately equal number of instances in the data set (i.e. have the number of instances in each class approximately n/c). We assume the time complexity of the base learning algorithm is linear in the number of training instances, and hence we can ignore the effect of the number of attributes.

In the worst case, a completely and/or balanced nested dichotomy can degenerate into a linear, and hence optimal, number of building a 1-class classifier based on half of the classes and the above analysis applies.

$$\begin{aligned}
 \sum_{i=0}^{c-2} \frac{c-i}{c} n &= \frac{n}{c} \sum_{i=0}^{c-2} c-i = \frac{n}{c} \left((c-1)c - \sum_{i=0}^{c-2} i \right) = \frac{n}{c} \left((c-1)c - \frac{(c-2)(c-1)}{2} \right) \\
 &> \frac{n}{c} \left((c-1)c - \frac{c(c-1)}{2} \right) = \frac{(c-1)}{2} n.
 \end{aligned}$$

Hence the worst-case time complexity is linear in the number of instances and classes.

Let us now consider the balanced case. Assuming c is even, we have $\log c$ layers of internal nodes. In each layer, all the training data need to be processed

(because the union of all b ’s in each layer is the original data a). Given h , we have a sorted h and the balanced leaf node is linearly in the number of instances, the overall n is become $n \log c$, i.e. is logarithmic in the number of classes and linear in the number of instances.

Adding a balancing algorithm whose time complexity would be linear, the advantage of the balanced tree becomes even more pronounced (because the size of the b ’s of data considered at each node decreases exponentially in height). No external sorting or even distribution of classes is not strictly necessary. This can be seen by considering the worst case, where one class has all of the instances. The worst case is the complexity of the unbalanced case, plain linear in the number of classes in height, and the one for the balanced case, logarithmic in the number of classes.

However, in the case of a sorted class distribution it is possible to improve on the class-balanced tree when the balancing algorithm is not linear. In this case it is a c -ary tree to divide the number of instances as evenly as possible at each node, so a sorted case is a c -ary tree of data considered at a node is also possible. This is why we have investigated the data-balanced approach discussed above.

2.4 Caching Models

The efficiency of our algorithm over the training is the feature of nested dichotomy. It is efficient for the fact that the leaf nodes are c -ary. Consider Fig. 1. In both cases, a class is highlighted in leaf nodes 2 and 5. The class will be identical because they are based on exactly the same data. It is not possible to build the tree. Consequently we can cache model that have been built in a hash table and used for a similar collection of features.

As explained by Frank and Kape [1], the size is $(3^c - (2^{c+1} - 1))/2$ possible c -ary nodes for a c -ary node, i.e. grows exponentially in the number of classes. Hence we can expect caching only to be a difference for relatively small number of classes. If we consider balanced dichotomy, the number of possible nodes is reduced. Consequently caching will be more beneficial in the balanced case.

3 Experiments

In the following we empirically investigate the effect of our proposed modification on accuracy. We used 21 multi-class UCI datasets [5]. The number of classes varies from 3 to 26. We also explored our extension with a artificial data that exhibits a large number of classes. For each case, we used 10 ensemble members (i.e. 10 copies of nested dichotomy are generated and their probability is averaged out for a prediction). J48, the implementation of the C4.5 decision tree learner [6] from the Weka toolbox [7] was used as the baseline for the extension.

Table 2. Effect of model caching on ENDS for UCI datasets

Dataset	Number of classes	Number of instances	Training time for ENDS	
			w/o caching	with caching
iris	3	150	0.03 ± 0.02	0.01 ± 0.00
balance-scale	3	625	0.28 ± 0.06	0.09 ± 0.04 ●
splice	3	3190	4.56 ± 0.32	1.37 ± 0.14 ●
waveform	3	5000	38.78 ± 0.65	11.42 ± 0.96 ●
lymphography	4	148	0.06 ± 0.02	0.04 ± 0.02 ●
vehicle	4	846	1.87 ± 0.11	1.08 ± 0.16 ●
hypothyroid	4	3772	3.13 ± 0.47	1.75 ± 0.29 ●
anneal	6	898	0.99 ± 0.08	0.82 ± 0.13 ●
zoo	7	101	0.07 ± 0.02	0.06 ± 0.02
autos	7	205	0.40 ± 0.05	0.36 ± 0.05
glass	7	214	0.30 ± 0.03	0.27 ± 0.03
segment	7	2310	6.61 ± 0.27	5.87 ± 0.37 ●
ecoli	8	336	0.25 ± 0.04	0.23 ± 0.04
optdigits	10	5620	72.53 ± 3.30	68.70 ± 3.00 ●
pendigits	10	10992	49.30 ± 2.00	47.07 ± 2.12 ●
vowel	11	990	4.21 ± 0.11	4.04 ± 0.16 ●
arrhythmia	16	452	21.14 ± 1.01	20.76 ± 1.09
soybean	19	683	1.02 ± 0.07	0.99 ± 0.06
primary-tumor	22	339	0.63 ± 0.06	0.63 ± 0.06
audiology	24	226	0.74 ± 0.06	0.74 ± 0.05
letter	26	20000	317.53 ± 11.44	315.74 ± 11.53

All experiments are performed for 10 runs of a fixed 5-fold cross-validation (UCI datasets) or 3-fold cross-validation (academic datasets). We also report standard deviation for the 50 (UCI datasets) or 30 (academic datasets) individual experiments. Results were obtained on a machine with a Pentium 4 3 GHz processor, running the Java HotSpot Client VM (build 1.4.2_03) on Linux, and implemented in second order. We observed significant differences using the corrected Levene test [8].

3.1 Applying Caching to ENDS

In this section we discuss the effect of caching individual classifiers in an ensemble of nested dichotomies. Table 2 has the average training time for ENDs with and without caching based on a hash table. Significant improvements in training time obtained by caching are marked with a ●.

The results show that the number of significant improvements for 14 of the 21 UCI datasets (of course, accuracy remains identical). The improvements are especially obvious for datasets with a large number of classes and a large number of instances. With a large number of classes one is more likely to encounter the same binary classifier in different parts of nested dichotomies in the ensemble. With a large number of instances, the problem is alleviated by avoiding rebuilding the same binary classifier. For instance, the training time on the waveform dataset, which

Table 3. Comparison of training time on UCI datasets

Dataset	Number of classes	Training time		
		ENDs	ECBNDs	EDBNDs
iris	3	0.01 ± 0.00	0.03 ± 0.02	0.02 ± 0.00
balance-scale	3	0.09 ± 0.04	0.09 ± 0.04	0.09 ± 0.04
splice	3	1.37 ± 0.14	1.45 ± 0.12	1.33 ± 0.20
waveform	3	11.42 ± 0.96	11.31 ± 0.56	11.24 ± 1.10
lymphography	4	0.04 ± 0.02	0.03 ± 0.02	0.03 ± 0.00
vehicle	4	1.08 ± 0.16	0.51 ± 0.06	0.52 ± 0.05
hypothyroid	4	1.75 ± 0.29	0.86 ± 0.12	0.92 ± 0.22
anneal	6	0.82 ± 0.13	0.63 ± 0.09	0.50 ± 0.13
zoo	7	0.06 ± 0.02	0.06 ± 0.03	0.06 ± 0.02
autos	7	0.36 ± 0.05	0.26 ± 0.04	0.25 ± 0.05
glass	7	0.27 ± 0.03	0.21 ± 0.04	0.20 ± 0.04
segment	7	5.87 ± 0.37	4.88 ± 0.34	4.98 ± 0.41
ecoli	8	0.23 ± 0.04	0.20 ± 0.03	0.21 ± 0.04
optdigits	10	68.70 ± 3.00	55.17 ± 1.91	55.03 ± 2.16
pendigits	10	47.07 ± 2.12	37.95 ± 1.53	38.40 ± 1.52
vowel	11	4.04 ± 0.16	3.62 ± 0.11	3.70 ± 0.12
arrhythmia	16	20.76 ± 1.09	19.20 ± 1.08	17.39 ± 1.56
soybean	19	0.99 ± 0.06	0.87 ± 0.06	0.85 ± 0.07
primary-tumor	22	0.63 ± 0.06	0.54 ± 0.06	0.54 ± 0.06
audiology	24	0.74 ± 0.05	0.64 ± 0.08	0.63 ± 0.09
letter	26	315.74 ± 11.53	273.45 ± 16.75	274.07 ± 16.84

has 5000 instances and only 3 classes, decreased dramatically for 38.78 seconds to 11.42 seconds by using a hashable. On the other hand, for the a_hypothyroid dataset, which has 16 classes and only 452 instances, the training time decreased only slightly, from 21.14 seconds to 20.76 seconds. From Table 2, we also see that the efficiency can improve even for the training time when the number of classes exceed 11. The chance of encoding the a binary class in a hashable becomes limited as the number of possible words grows. Moreover, the number of instances in the dataset (excluding the le_data) is still on the order of magnitude by using a hashable is not noticeable. We also see, for example, in the case of the a class and the a class are essentially no difference in the number of instances.

3.2 Comparing ENDs, ECBNDs, and EDBNDs

As we have seen, caching does not help when the a any class. In the following we will see how a balanced nested dichotomy helps in the case. We will also look at training time and then the effect on accuracy.

Training time. Table 3 shows the training time for END, class-balanced END (ECBND), and data-balanced END (EDBND), on the UCI datasets. Model caching was applied in all the evaluation of END. A • indicates a significant reduction in the number of END.

Table 4. Comparison of training time on artificial datasets

Number of classes	Number of instances	Training time		
		ENDs	ECBNDs	EDBNDs
10	820	0.60 ± 0.09	0.58 ± 0.07	0.58 ± 0.07
20	1390	1.50 ± 0.12	1.42 ± 0.08	1.44 ± 0.09
30	1950	2.72 ± 0.11	2.31 ± 0.12	2.33 ± 0.12
40	2410	3.87 ± 0.16	3.18 ± 0.14	3.24 ± 0.13
50	3090	5.55 ± 0.23	4.54 ± 0.17	4.57 ± 0.20
60	3660	7.48 ± 0.29	5.86 ± 0.17	5.90 ± 0.25
70	4560	10.41 ± 0.36	8.23 ± 0.28	8.35 ± 0.30
80	5010	12.32 ± 0.43	9.56 ± 0.33	9.67 ± 0.31
90	5840	15.75 ± 0.53	12.62 ± 0.44	12.78 ± 0.34
100	6230	18.25 ± 0.61	13.61 ± 0.38	13.98 ± 0.44
150	9590	40.65 ± 1.90	27.63 ± 0.77	28.19 ± 0.65
200	12320	66.41 ± 2.95	42.37 ± 1.30	42.70 ± 1.30

The results show that using class-balanced nested dichotomies is significantly reduced training time on 14 of the 21 datasets. Using the data-balanced checker also helped: EDBND is significantly more efficient than END on 14 datasets, just like ECBND. Compared to class-balanced trees, data-balanced trees are significantly more efficient on one dataset (a, h, h, ia). (This information is not included in Table 3.) This dataset has an extremely unbalanced class distribution and this is why the data-balanced approach helped.

The advantage of the balanced checker is exercised on datasets with more than 3 classes. On high-class datasets, all nested dichotomies are class-balanced, so we would not expect any significant difference between END and ECBND. The experimental results bear this out.

Table 4 shows the training time for 12 artificial datasets. To generate the datasets we used a class generator and varied the number of classes from 10 to 200. In advance in the algorithm we assigned the class label. Each instance in the dataset consists of one boolean attribute and won't be a class. The attributes value range will be different but could overlap. An attribute value will be generated and only within each class. The number of instances in each class (i.e. class) was also generated and varied between 20 and 110.

The results on the artificial datasets show that the balanced checker exhibits a significant advantage in terms of training time when 30 or more classes are present in the data. There was no significant difference in training time for the two balanced checkers (class-balanced vs. data-balanced) on any of the datasets. This indicates that the class distribution in our artificial datasets is not skewed enough for the data-balanced approach to help.

Accuracy. I am interested in finding out if they affect accuracy in a significant fashion. Hence I did a one-way ANOVA to evaluate the effect of nested dichotomies on accuracy. Table 5 shows the estimated accuracy for END, ECBND, and EDBND on the UCI datasets. We can see that there is no data set with a significant difference in accuracy for END and ECBND. This is the

Table 5. Comparison of accuracy on UCI datasets

Dataset	Number of classes	Percent correct		
		ENDs	ECBNDs	EDBNDs
iris	3	94.13 ± 3.84	94.13 ± 3.72	94.27 ± 3.81
balance-scale	3	79.92 ± 2.37	79.49 ± 2.41	79.78 ± 2.31
splice	3	94.75 ± 1.01	94.55 ± 0.98	93.07 ± 1.33
waveform	3	77.89 ± 1.88	77.53 ± 1.91	77.85 ± 2.06
lymphography	4	77.73 ± 7.47	76.63 ± 6.35	76.90 ± 6.93
vehicle	4	73.20 ± 2.92	72.36 ± 2.30	72.36 ± 2.30
hypothyroid	4	99.54 ± 0.26	99.51 ± 0.27	99.54 ± 0.28
anneal	6	98.63 ± 0.80	98.44 ± 0.75	98.53 ± 0.62
zoo	7	93.66 ± 5.67	93.87 ± 4.61	93.88 ± 4.50
autos	7	76.20 ± 6.11	74.83 ± 6.62	75.32 ± 7.10
glass	7	72.82 ± 7.42	73.51 ± 6.17	72.25 ± 6.84
segment	7	97.45 ± 0.83	97.35 ± 0.80	97.39 ± 0.87
ecoli	8	85.60 ± 4.11	85.36 ± 4.06	84.88 ± 4.13
optdigits	10	96.99 ± 0.49	97.14 ± 0.45	97.18 ± 0.50
pendigits	10	98.59 ± 0.27	98.76 ± 0.25	98.76 ± 0.26
vowel	11	88.31 ± 2.66	89.98 ± 2.47	89.24 ± 2.79
arrhythmia	16	72.59 ± 3.24	72.82 ± 4.11	71.51 ± 3.55
soybean	19	93.90 ± 1.63	94.49 ± 1.69	94.36 ± 1.78
primary-tumor	22	44.72 ± 5.04	46.28 ± 4.61	45.96 ± 4.62
audiology	24	78.46 ± 5.44	79.66 ± 5.12	79.48 ± 5.23
letter	26	94.33 ± 0.37	94.50 ± 0.36	94.51 ± 0.35

Table 6. Comparison of accuracy on artificial datasets

Number of classes	Number of instances	Percent correct		
		ENDs	ECBNDs	EDBNDs
10	820	78.08 ± 1.94	78.34 ± 2.35	78.32 ± 2.32
20	1390	77.79 ± 1.87	77.21 ± 1.44	77.47 ± 1.66
30	1950	77.09 ± 1.61	76.93 ± 1.53	76.85 ± 1.46
40	2410	76.64 ± 1.24	76.56 ± 1.39	76.46 ± 1.24
50	3090	76.26 ± 1.09	76.17 ± 1.26	76.25 ± 1.19
60	3660	76.43 ± 1.08	76.33 ± 1.04	76.37 ± 0.95
70	4560	73.58 ± 1.12	73.27 ± 0.97	73.50 ± 0.90
80	5010	75.85 ± 1.06	75.61 ± 0.94	75.71 ± 0.87
90	5840	76.41 ± 0.84	76.40 ± 0.91	76.41 ± 0.87
100	6230	76.59 ± 0.77	76.54 ± 0.73	76.50 ± 0.85
150	9590	75.92 ± 0.66	75.89 ± 0.72	75.86 ± 0.62
200	12320	75.89 ± 0.51	75.67 ± 0.51	75.73 ± 0.49

de i, ed o co e. Fo EDBND , he, e i one h ee-cla da a e (lice) whe e he acc \acy i igni can ly ,ed ced co a, ed o END . The lice da a ha a ewed cla di ,ib ion, whe e one cla ha abo half he in ance and he e i evenly di ,ib ed a ong he ,e aining wo cla e. We ea ,ed he dive, i y of he h, ee y e of en e ble on hi da a e ,ing he a a

is a i c. This a i c can be ed o ea e agree en be ween ai. of en e ble e be. [9]. Fo EDBND , he ean a a val e ove all ai. , ead on he aining da a, wa 0.96, which wa indeed highe han he ean a a val e fo END and ECBND (0.94 and 0.93 e ec ively). Thi indica e ha ed c ion in dive. i y i he ead on fo he d o in e fo. ance.

Table 6. how he a e info. a ion fo he a i cial da a e. In hi ca e he e i no a i ngle da a e whe e he e i a i gni can diffe nce in acc. acy be ween any of he che e.

4 Conclusions

En e ble of ne ed dicho o ie have e en ly been hown o be a ve, y o i ing e a lea ning che e fo. l i-cla. o ble. . They od ce acc. a e cla. i ca ion and yield cla. obabili e e i a e in a na. al way. In hi a e, we have hown ha i i o ble o i o ve he n i e of hi e a lea ning che e wi ho affec ing acc. acy. A i le way o i o ve n i e fo o ble. wi ha. all n. be. of cla e i o cache wo-cla. o del and e e he in diffe n. e be. of an en e ble of ne ed dicho o ie. On o ble. wi h any cla e we have hown ha i n g cla.-balanced ne ed dicho o ie i gni can ly i o ve n i e, wi no i gni can change in acc. acy. We have al o e en ed a da a-balanced che e ha can hel o i o ve n i e f, he, when he e a e. any cla e and he cla. di. ib ion i highly e wed.

References

1. Frank, E., Kramer, S.: Ensembles of nested dichotomies for multi-class problems. In: Proc Int Conf on Machine Learning, ACM Press (2004) 305–312
2. Dietterich, T., Bakiri, G.: Solving multiclass learning problems via error-correcting output codes. *Journal of Artificial Intelligence Research* **2** (1995) 263–286
3. Fürnkranz, J.: Round robin classification. *Journal of Machine Learning Research* **2** (2002) 721–747
4. Fox, J.: *Applied Regression Analysis, Linear Models, and Related Methods*. Sage (1997)
5. Blake, C., Merz, C.: UCI repository of machine learning databases. University of California, Irvine, Dept. of Inf. and Computer Science (1998) [www.ics.uci.edu/~mllearn/MLRepository.html].
6. Quinlan, J.: *C4.5: Programs for Machine Learning*. Morgan Kaufmann, Los Altos, CA (1992)
7. Witten, I.H., Frank, E.: *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann (2000)
8. Nadeau, C., Bengio, Y.: Inference for the generalization error. *Machine Learning* **52** (2003) 239–281
9. Dietterich, T.G.: An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. *Machine Learning* **40** (1998) 139–157

Protein Sequence Pattern Mining with Constraints

Pedro Gabriel Ferreira* and Paulo J. Azevedo**

University of Minho, Department of Informatics,
Campus of Gualtar, 4710-057 Braga, Portugal
{pedrogabriel, pja}@di.uminho.pt

Abstract. Considering the characteristics of biological sequence databases, which typically have a small alphabet, a very long length and a relative small size (several hundreds of sequences), we propose a new sequence mining algorithm (*gIL*). *gIL* was developed for linear sequence pattern mining and results from the combination of some of the most efficient techniques used in sequence and itemset mining. The algorithm exhibits a high adaptability, yielding a smooth and direct introduction of various types of features into the mining process, namely the extraction of rigid and arbitrary gap patterns. Both breadth or a depth first traversal are possible. The experimental evaluation, in synthetic and real life protein databases, has shown that our algorithm has superior performance to state-of-the art algorithms. The use of constraints has also proved to be a very useful tool to specify user interesting patterns.

1 Introduction

In the development of sequence pattern mining algorithms, two communities can be considered: the *Data Mining* community and the *Bioinformatics* community. The algorithms from the Data Mining community inherited some characteristics from the association rule mining algorithms. They are best suited for data with many (from hundred of thousands to millions) sequences with a relative small length (from 10 to 20), and an alphabet of thousands of events, e.g. [9,7,11,1]. In the bioinformatics community, algorithms are developed in order to be very efficient when mining a small number of sequences (in the order of hundreds) with large lengths (few hundreds). The alphabet size is typically very small (ex: 4 for DNA and 20 for protein sequences). We emphasize the algorithm Teiresias [6] as a standard.

The major problem with Sequence pattern mining is that it usually generates too many patterns. When databases attain considerable size or when the average

* Supported by a PhD Scholarship (SFRH/BD/13462/2003) from Fundação Ciência e Tecnologia.

** Supported by Fundação Ciência e Tecnologia - Programa de Financiamento Pluri-anual de Unidades de I & D, Centro de Ciências e Tecnologias da Computação - Universidade do Minho.

length of the sequences is very long, the mining process becomes computationally expensive or simply infeasible. This is often the case when we are mining biological data like proteins or DNA. Additionally, the user interpretation of the results turns out to be a very hard task since the interesting patterns are blurred into the huge amount of outputted patterns. The solution to this problem can be achieved through the definition of alternative interesting measures besides support, or with user imposed restrictions to the search space. When properly integrated in the mining process these restrictions reduce the computation demands in terms of time and memory, allowing to deal with datasets that are otherwise potentially untractable. These restrictions are expressed through what is typically called as constraints. The use of Constraints enhances the database queries. The runtime reduction grants the user with the opportunity to interactively refine the query specification. This can be done until an expected answer is found.

2 Preliminaries

We consider the special case of linear sequences databases. A database D is as a collection of linear sequences. A sequence $S = \langle e_1 e_2 \dots e_n \rangle$ is a sequence composed by successive atomic elements, generically called events. Examples of this type of databases are protein or DNA sequences or website navigation paths. The term *linear sequence* is used to make the distinction from the transactional sequences, that consist in sequences of events (usually called as *transaction*). Given a sequence S , S' is subsequence of S if S' can be obtained by deleting some of the events in S . A sequence pattern is called a *linear sequence pattern* if it is found to be subsequence of a number of sequences in the dataset greater or equal to a specified threshold value. This value is called, *support threshold*, σ , and is defined as an user parameter. The ids represents the list of sequence identifiers where the pattern occurs. The cardinality of this list corresponds to the support of that pattern.

Considering patterns in the form $A_1 - x(p_1, q_1) - A_2 - x(p_2, q_2) - \dots - A_n$, a sequence pattern is an *event-gap-event-gap-event* when a variable (zero or more) number of gaps exist between adjacent events in the pattern, i.e. $p_i \leq q_i, \forall i$. Typically a variable gap with n minimum and m maximum number of gaps is described as $-x(n, m)-$. In the sequences $\langle 1\ 5\ 3\ 4\ 5 \rangle$ and $\langle 1\ 2\ 2\ 3 \rangle$ exists an arbitrary gap pattern $1-x(1, 2)-3$. A *rigid gap pattern* is a pattern where gaps contain a fixed size for all the database occurrences of the sequence pattern, i.e. $p_i = q_i, \forall i$. To denote a rigid gap the $-r(n)-$ notation is used, where n is the size of the gap. The $1-r(2)-3$ is a pattern of length 4, in the sequences $\langle 1\ 2\ 5\ 3\ 4\ 5 \rangle$ and $\langle 1\ 1\ 6\ 3 \rangle$. Each gap position is denoted by the "." (wildcard) symbol, meaning that it matches any symbol of the alphabet. A pattern belongs to one of three classes: *primitive*, *extension* or *maximal*. A sequence pattern is *primitive* if it is not contained in any other pattern, and *extension* when all its extensions have an inferior support than itself. The *maximal* refers to when all the patterns are enumerated. When extending a sequence pattern $S = \langle s_1\ s_2 \dots s_n \rangle$, with a new event s_{n+1} , then

S is called a *sequence pattern* and $S' = \langle s_1 s_2 \dots s_n s_{n+1} \rangle$ the *sequence*. If an event b occurs after a in a certain sequence, we denoted it as: $a \rightarrow b$, and a is called the *predecessor*, $pred(a \rightarrow b) = a$, and b the *successor*, $succ(a \rightarrow b) = b$. The pair is frequent if it occurs in at least σ sequences of the database.

Constraints represent an efficient way to prune the search space [9,10]. Considering the user's point of view, it also enables to focus the search on more interesting sequence patterns. The most common and generic types of constraints are:

- *Event set constraint*: restricts the set of the events ($\{e_1, e_2, \dots, e_n\}$) that may appear in the sequence patterns,
- *Distance constraint*: defines the (min_dist, max_dist) minimum distance or the maximum distance (min_dist, max_dist) that may occur between two adjacent events in the sequence patterns,
- *Span constraint*: defines the maximum distance (max_span) between the first and the last event of the sequence patterns.
- *Start constraint*: determines that the extracted patterns start with the specified events ($\{e_1, e_2, \dots, e_n\}$).

Another useful feature in sequence mining, in particular to protein pattern mining, is the use of *substitution sets*. When used during the mining process an event can be substituted by another event belonging to the same set. A *hierarchy* of relations can be represented through substitution sets.

Depending on the target application of the frequent sequence patterns other measures of interest and scoring can be applied as posterior step of the mining process. Since the closed and the maximal patterns are not necessarily the most interesting we designed our algorithm in order to find all the frequent patterns. From the biological point of view, rigid patterns allow to find more well conserved regions, while arbitrary patterns permit the cover of a large number of sequences in the database.

The problem we address in this paper can be formulated as follow: given a database D of linear sequences, a minimum support, σ , and the optional parameters: *event set constraint*, *distance constraint*, *span constraint* and *start constraint*, find all the *closed* or *maximal*, frequent sequence patterns that respect the defined constraints.

3 Algorithm

The proposed algorithm uses a Bottom-Up search space enumeration and a combination of frequent pairs of events to extend and find all the frequent sequence patterns. The algorithm is divided in two phases: *enumeration* and *extension*. Since the frequent sequences are obtained from the set of frequent pairs, the first phase of the algorithm consists in traversing all the sequences in the database and building two auxiliary data structures. The first structure contains the set of all pairs of events found in the database. Each pair representation points to the sequences where they appear (through a sequence identifier bitmap). The second data structure consists of a vertical representation

of the database. It contains the positions or offsets of the events in the sequences where they occur. This information is required to ensure that the order of the events along the data sequence is respected. Both data structures are thought for quick information retrieval. At the end of the scanning phase we obtain a map of all the pairs of events present in the database and a vertical format representation of the original database. In the second phase the pairs of events are successively combined to find all the sequence patterns. These operations are fundamentally based on two properties:

Property 1 (Anti-Monotonic).

Property 2 (Sequence Transitive Extension).

$S = \langle s_1 \dots s_n \rangle, C_S$
 O_S
 $C_S \dots P = (s_j \rightarrow s_m), C_P \dots O_P \dots succ(P)$
 $C_P \dots succ(S) = pred(P), \dots s_n = s_j, \dots$
 $E = \langle s_1 \dots s_n s_m \rangle, \dots C_E, \dots C_E = \{X : \forall X \text{ in } C_S \cap C_P, O_P(X) > O_S(X)\}$

Hence, the basic idea is to successively extend a frequent pair of events with another frequent pair, as long as the predecessor of one pair is equal to the successor of the other. This joining step is sound provided that the above mentioned properties (1 and 2) are respected. The joining of pairs combined with a breadth first or a depth first traversal yields all the frequent sequences patterns in the database.

3.1 Scanning Phase

The first phase of the algorithm consists in the following procedure: For each sequence in D , obtain all ordered pairs of events, without repetitions. Consider the sequence 5 in the example database of table 1(a). The obtained pairs are: $1 \rightarrow 2, 1 \rightarrow 3, 1 \rightarrow 4, 2 \rightarrow 2, 2 \rightarrow 3, 2 \rightarrow 4$ and $3 \rightarrow 4$. During the determination of the pairs of events the first auxiliary data structure, that consists of an N-bidimensional matrix, is built and updated. N corresponds to the size of the alphabet. The N^2 cells in the matrix correspond to the N^2 possible combinations of pairs. We call this structure the $Cell(i, j)$. Each $Cell(i, j)$ contains the information relative to the pair $i \rightarrow j$. This information consists of a bitmap that indicates the presence (1) or the absence (0) in the respective sequence

Table 1. (a) Parameters used in the synthetic data generator; (b) Properties of the proteins datasets

Symbol	Meaning	DataSet	NumSeq	AlphabetSize	AvgLen	MinLen	MaxLen
S	Number of Sequences (x 10^3)	Yeast	393	21	256	15	1859
L	Avg. Length of the sequences	PSSP	396	22	158	21	577
R	Alphabet Size	nonOM	60	20	349	53	1161
P	Distribution Skewness	mushroom	8124	120	23	23	23

(i -th bit corresponds to the sequence i in D) and an integer that contains the support count. This last value allows a fast support checking. For each pair $i \rightarrow j$ we update the respective $Cell(i, j)$ in the Bitmap Matrix, by activating the bit corresponding to the sequence where the pair occurs and incrementing the support counter. As an example, for the pair $1 \rightarrow 3$, the $Cell(1, 3)$ is represented in figure 1(b):

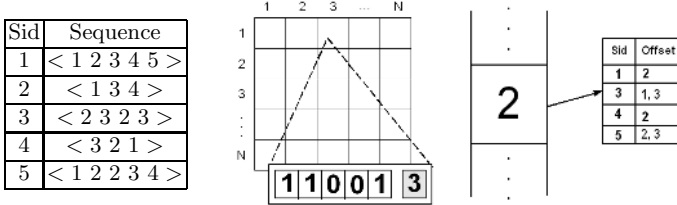


Fig. 1. (a) An example database; (b) Content of the Cell(1,3) in the Bitmap Matrix; (c) Representation of event 2 in the Offset Matrix

This means that the pair occurs in the database sequence 1, 2 and 5 and has a support of 3. Simultaneously, as each event in the database is being scanned, a second data structure called \dots is also being built. Conceptually, this data structure consists of an adjacency matrix that will contain all the offset (positions) of all the events in the entire database. Each event is a key that points to a list of pairs $\langle \dots, \dots \rangle$, where \dots is a list of all the positions of the event in the sequence \dots . Thus, the Offset Matrix is a vertical representation of the database. Figure 1(c) shows the information stored in the Offset Matrix for the event 2.

3.2 Sequence Extension Phase

We start by presenting how arbitrary gap patterns are extracted. In section 3.4 we will show how easily our algorithm can be adapted to extract rigid gap patterns. For implementing the extension phase we present two tests (algorithms) that conjunctively are necessary and sufficient conditions to consider as frequent a new extended sequence.

This is a quick test that implements property 1. The \dots function gets the correspondent bitmaps of S and P . The intersection operation is also very fast and simple and the support function retrieves the support of the intersection bitmap. This test allows the verification of a \dots condition for the extended sequence to be frequent. A second test is necessary

```

input : ( ... ); ( ... ); ( ... )
= ( ... ) and = ( ... )
if ( ... ) ≥ , then return OK.

```

Algorithm 1: Support Test

```

input : r( ); ( ); ( ); ( )
1 = ( r);
2 = ( );
3 = 0;
4 foreach Sid in seqLst do
5   = offsetLst( , );
6   = offsetLastEvent( , );
7   = offsetStartEvent( , );
8   if ∃ ∈ = , +1; then
9     if ( - ) then = ( - )
10    if ( - ) then = ( - )
11  end
12    ( , );
13 end
14 if ≥ then
15   return OK;
16 end

```

Algorithm 2: Order Test

to guarantee that the order of the events is kept along the sequences that $C_{S'}$ bitmap points to.

Algorithm 2 assumes that, for each frequent sequence, additional information besides the sequence event list is kept during the extension phase. Namely, the corresponding bitmap that for the case exposed in algorithm 1 will be $C_{S'}$ if S' is determined to be frequent. Also two offset lists in the form $\langle \dots \rangle$ are kept. One will contain the offset of the last event of the sequence, \dots , and will be used for the "Order Test". The second, \dots , contains the offset of the first event of the sequence pattern in all the Sid where it appears. This will be used when the verification of the window constraint is performed. In the Order Test, given a bitmap resulted from the support test, the \dots function returns the list of the sequence identifiers for the bitmap. The function \dots returns a list of offset values of the event in the respective Sid. For each sequence identifier it is tested whether the extension pair has an offset greater than the offset value of the extended sequence. This implements the computation of C_E and the $offsetList$ of $succ(E)$ as in property 2. At line 13 the $diffTest$ function performs a simple test to check whether the minimum support is still reachable. At the end of the procedure (lines 15 to 17) it is tested whether the order of the extended sequence pattern is respected in a sufficient number of database sequences. In the positive case the extended sequence is considered frequent. Given algorithm 1 and 2, property 3 guarantees the necessary and sufficient conditions to safely extend a base sequence into a frequent one.

Property 3 (Frequent Extended Sequence). $\dots \sigma$, $S = \langle E_1 \dots E_n \rangle$, $|S| \geq 2$, $P = E_k \rightarrow E_w \dots E_n = E_k \dots S' = \langle E_1 \dots E_n g_{n,k} E_k \rangle$, $g_{n,k} = -x(n, m) - r(n) - \dots OK$

3.3 Space Search Traversal

Guided by the Bitmap Matrix the search space can be traversed using two possible approaches: \dots or \dots . For both cases the set of the frequent

sequences starts as the set of frequent pairs. In the depth first mode it starts with a sequence of size 2 that is successively expand until it can not be further extended. Then we backtrack and start extending another sequence. The advantage of this type of traversal is that we don't need to keep all the intermediary frequent sequence information, in contrast with the breadth first traversal where all the information of the sequences size k need to be kept before the sequences of size $k+1$ are generated. This yields in some cases, a significant memory reduction.

3.4 Rigid Gap Patterns

The algorithm described in 2 is designed to mine arbitrary gap patterns. Using *gIL* to mine rigid gap patterns requires only minor changes in the Order Test algorithm. Lines 4 to 11 in algorithm 2 are rewritten in algorithm 3. In this algorithm, first it is collected (in *gapLst*) the size of all the gaps for a certain sequence extension. Next, for each gap size it is tested whether the extended sequence is frequent. One should note that for rigid gap patterns, two sequence patterns with the same events are considered different if the gaps between the events have different size, e.g., $\langle 1 \cdot \cdot 2 \rangle$ is different from $\langle 1 \cdot \cdot \cdot 2 \rangle$.

```

1  foreach Sid in seqLst do
2      = offsetLst( , );
3      = offsetLastEvent( , );
4      = offsetStartEvent( , );
5      if  $\exists \in$  , - ; then
6          end
7      end
8  end
9  foreach R in gapLst do
10     foreach Sid in seqLst do
11         Repeat Step 2 to 4;
12         if  $\exists \in$  , ( - ) = then
13             end
14         end
15     end
16     if  $\geq$  then
17         return OK;
18     end
19 end

```

Algorithm 3: Algorithm changes to mine rigid gap patterns

4 Constraints

The introduction of constraints in the *gIL* algorithm like $\langle 1 \cdot \cdot \cdot 2 \rangle$, $\langle 1 \cdot \cdot \cdot 2 \rangle$, $\langle 1 \cdot \cdot \cdot 2 \rangle$ is a straightforward process and translates into a considerable performance gain. These efficiency improvements are naturally expected since (depending on the values of the constraints) the search space can be greatly reduced. The introduction of $\langle 1 \cdot \cdot \cdot 2 \rangle$ is also very easy to achieve. Implementing events exclusion constraint and substitution sets turns out to be a natural operation. Simple changes in the Bitmap Matrix (that guides the sequence extension) and in the Offset Matrix (discriminates the positions of the events in every sequence where they occur) enable this implementations. The new features are introduced between the scanning phase and the sequence extension phase. The min/max gap and window constraints constitute an additional to be applied when the sequence is extended.

4.1 Events Exclusion, Start Events and Substitution Sets

The event exclusion constraint is applied by traversing the rows and columns of the Bitmap Matrix where the excluded events occurs. At that positions the support¹ count variable in the respective cells is set to zero. Start events constraints are also straightforwardly implemented by allowing extensions only to the events in *StartEventSets*.

When substitution sets are activated, one or more sets of equivalent events are available. For each set of equivalent events one has to form the union of the rows (horizontal union) and columns (vertical union) in the Bitmap Matrix, where those events occur. The vertical union is similar to the horizontal union. Moreover, for all the equivalent events, one needs to pairwise intersect the sequences where they occur and then perform the union of the offsetLists for the intersected sequences. This results in the new offsetLists of the equivalent events.

4.2 Min / Max Gap and Window Size

These constraints are trivially introduced in the "Order Test". In algorithm 2, the test in line 8 is extended with three additional tests: $(X - Y) < maxGap$ AND $(X - Y) > minGap$ AND $(X - W) < windowSize$.

5 Experimental Evaluation

We evaluated our algorithm along different variables using two collections of synthetic and real datasets. To generate the synthetic datasets we developed a sequence generator based on the Zipfian distribution. This generator receives the following parameters (see table 1(a)): number of sequences, average length of the sequences, alphabet size and a parameter p that expresses the skewness of the distribution. This generator has allowed us to generate sequences with a relative small alphabet. The evaluated variables for this datasets were: n , l , a , p , n , l , a , p , n , l , a , p , n , l , a , p . Additionally, we tested the mushroom dataset used at the FIMI workshop [4]. To represent real life data, we used several datasets of proteins. The Yeast (n , l , a , p) dataset is available at [5] and PSSP used for protein secondary structure prediction [3]. We also used a subset of the non Outer Membrane proteins obtained from [8]. The properties for this datasets are summarized in table 1(b). It is interesting to notice that, for all datasets, gIL's scanning phase time is residual (less than 0.4 seconds).

Since gIL finds two types of patterns we performed evaluation against two different algorithms. Both are in memory algorithms, assuming that the database completely fits into main memory. For the arbitrary gap patterns from the all patterns class we compared gIL with the SPAM [1] algorithm. SPAM has shown to outperform SPADE [11] and PrefixSpan [7] and is a state-of-the-art algorithm

¹ Future interactions on this dataset still have the Bitmap Matrix intact since the bitmaps remain unchanged.

in transactional sequence pattern mining. The datasets suffer a conversion into the transactional dataset format, in order to be processed by SPAM. In this conversion each customer is considered as a sequence and each contains a unique item (event).

For the rigid gap patterns we compared gIL with Teiresias [6], a well known algorithm from the bioinformatics community. It can be obtained from [2]. It is, as far as we know, the most complete and efficient algorithm for mining closed (called "most specific" in their paper) frequent rigid gap patterns. Closed patterns are a subset of all frequent sequence patterns. In this sense, gIL (which derives all patterns) tackles a more general problem and consequently considers a much larger search space than Teiresias. Besides minimum support, Teiresias uses two additional parameters. L and W are respectively the number of non-wild cards events in a pattern and the maximum spanning between two consecutive events. Since gIL starts by enumerating patterns with size 2, we will set $L=2$ and W to the maxGap value. All the experiments² were performed using exact discovery, i.e. without the use of substitution sets, and on a 1.5GHz Intel Centrino machine with 512MB of main memory, running windows XP Professional. The applications were written in C/C++ language.

5.1 Arbitrary Gap Patterns Evaluation

We start by comparing the efficiency of SPAM with the gIL algorithm without constraints. In figure 2(a) and 2(b) we tested different values of support for two datasets of $1K$ and $2K$ respectively. The sequences have an average length of 60 and an alphabet of 20 events. It was clear in these two experiments that for relative smaller dataset sizes and lower support values gIL becomes more efficient than SPAM. Figure 2(c) shows the scalability of the algorithms in respect to the dataset size for a support of 30%. This graphic shows that gIL scales well in relation to the dataset size.

In order to test a dataset with different characteristics, namely larger alphabet size, small length and greater dataset size, we used the Mushroom dataset, see figure 3(a). In figure 3(b) we have runtimes of gIL for datasets with one thousand sequences and different values of average sequence length. It was imposed a maxGap constraint of 15. As we observed during all the experiments, there is a critical point in the support variation, typically between 10% and 20%, that translates into an explosion of number of frequent patterns. This leads to an exponential behaviour in the algorithm's runtime. Even so, we can see that gIL shows similar behaviour for the different values of sequence length. Figure 3(c) measures the relative performance time, i.e. the ratio between the mining time with constraints and without constraints. These values were obtained for a support of 70%. Runtime without constraints was 305 seconds. It describes the behaviour of the algorithm when decreasing the maxGap and the Window values.

² Further details and results can be obtained from an extended version of this paper.

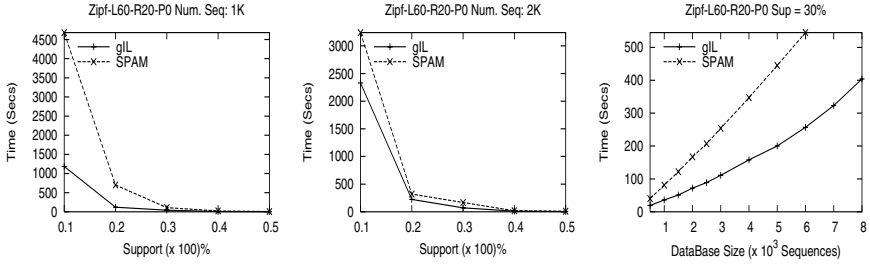


Fig. 2. (a) Support variation with Zipf database size=1K; (b) Support variation with Zipf database size=2K; (c) Scalability of gIL w.r.t. database size with a support of 30%

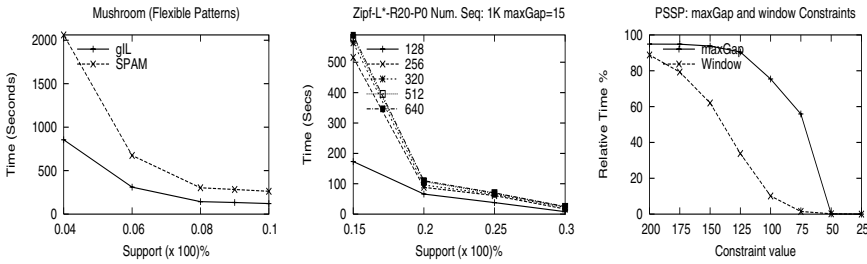


Fig. 3. (a) Support variation for the Mushroom dataset; (b) Scalability of gIL w.r.t sequence size for different support values (c) Performance evaluation using maxgap and windowgap constraints

In respect to memory usage both algorithms showed a low memory demand for all the datasets. For the Mushroom dataset which was the most demanding in terms of memory, SPAM used a maximum of 9 MB for a support of 4% and gIL a constant memory usage of 26 MB for all the support values. gIL shows a constant and support independent memory usage since once the data structures are built for a given dataset they remain unchanged.

5.2 Rigid Gap Patterns Evaluation

In order to assess the performance of gIL in the mining of rigid gap patterns we compared it with Teiresias [6], for different proteins datasets. In figure 4(a) and 4(b) the Yeast dataset was evaluated for two values of maxGap(W), 10 and 15. The results showed that gIL outperforms Teiresias by an order of magnitude. When comparing the performance of the algorithms in relation to the PSSP (figure 4(c)) and the nonOM (figure 5(a)) datasets, for a maxGap of 15, gIL outperforms Teiresias by a factor of 2 in the first case. This difference becomes more significant in the second case. The nonOM dataset has a greater average sequence length, but a small dataset size. This last characteristic results into a smaller bitmap length yielding a significant performance improvement. As

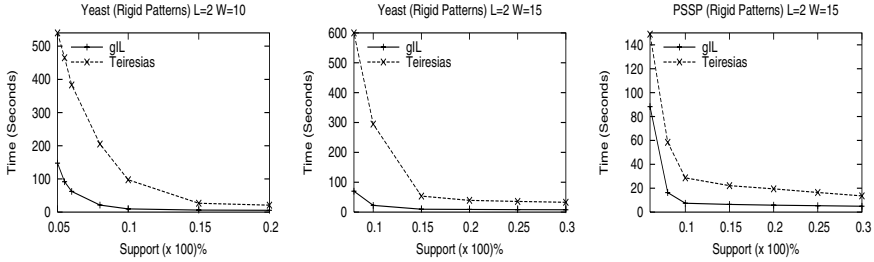


Fig. 4. (a) Support variation for the Yeast dataset, with $L=2$ and $W(\max\text{Gap}) = 10$; (b) Support variation for the Yeast dataset, with $L=2$ and $W(\max\text{Gap}) = 15$; (c) Support variation for the PSSP dataset, with $L=2$ and $W(\max\text{Gap}) = 15$

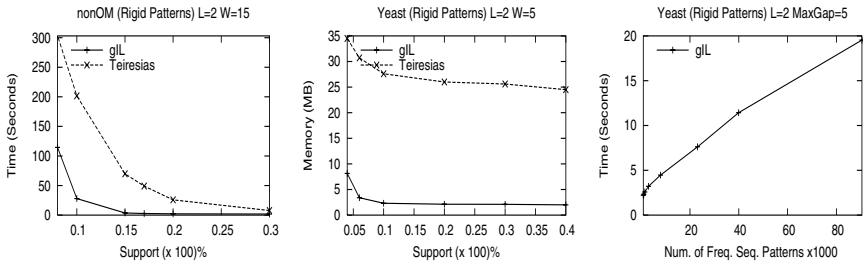


Fig. 5. (a) Support variation for the nonOM dataset, with $L=2$ and $W(\max\text{Gap}) = 15$; (b) Memory usage for the Yeast dataset, with $L=2$ and $W(\max\text{Gap}) = 5$; (c) Scalability of gIL w.r.t number of sequences for the Yeast dataset

we already verified in the arbitrary gap experiments, gIL memory usage maintains nearly constant for all the tested support values (figure 5(b)). Figure 5(c) shows the linear scalability of gIL in relation to the number of frequent sequence patterns.

6 Conclusions

We presented an algorithm called *gIL*, suitable to work with databases of linear sequences with a long average length and a relative small alphabet size. Our experiments showed that for the particular case of the proteins datasets, gIL exhibits superior performance to state-of-the-art algorithms. The algorithm has a high adaptability, and thus it was easily changed to extract two different types of patterns: arbitrary and rigid gap patterns. Furthermore, the data organization allows a straightforward implementation of constraints and substitution sets. These features are pushed directly into the mining process, which in some cases enables the mining in useful time of otherwise untractable problems. In this sense gIL is an interesting and powerful algorithm to be applied in a broader range of domains and in particular suitable for biological data. Thus, even when

performing extensions an event at a time (using a smart combination of some of the most efficient techniques that have been used in the task of itemset and sequence mining) one can obtain an algorithm that efficiently handles the explosive nature of pattern search, inherent to the biological sequence datasets.

References

1. J. Ayres, J. Flannick, J. Gehrke, and T. Yiu. Sequential pattern mining using a bitmap representation. In *Proceedings of the 8th SIGKDD International Conference on KDD and Data Mining, 2002*.
2. IBM Bioinformatics. Teiresias. <http://www.research.ibm.com/bioinformatics/>.
3. James Cuff and Geoffrey J. Barton. Evaluation and improvement of multiple sequence methods for protein secondary structure prediction. In *PROTEINS: Structure, Function, and Genetics*, number 34. WILEY-LISS, INC, 1999.
4. Fimi. Fimi workshop 2003 (mushroom dataset). <http://fimi.cs.helsinki.fi/fimi03>.
5. GenBank. yeast (*saccharomyces cerevisiae*). www.maths.uq.edu.au.
6. A. Floratos I. Rigoutsos. Combinatorial pattern discovery in biological sequences: the teiresias algorithm. *Bioinformatics*, 1(14), January 1998.
7. J. Pei, J. Han, B. Mortazavi-Asl, H. Pinto, Q. Chen, U. Dayal, and M.-C. Hsu. PrefixSpan: Mining sequential patterns efficiently by prefix projected pattern growth. In *Proceedings of the International Conference on Data Engineering, ICDE 2001*.
8. Psort. Psort database. <http://www.psort.org/>.
9. Ramakrishnan Srikant and Rakesh Agrawal. Mining sequential patterns: Generalizations and performance improvements. In *Proceedings 5th International Conference on Extending DataBase Technology, 1996*.
10. Mohammed J. Zaki. Sequence mining in categorical domains: Incorporating constraints. In *In Proceedings of 9th International Conference on Information and Knowledge Management, CIKM 2000*.
11. Mohammed J. Zaki. Spade: An efficient algorithm for mining frequent sequences. *Machine Learning*, 42(1-2):31-60, 2001.

An Adaptive Nearest Neighbor Classification Algorithm for Data Streams

Yan-Nei Law and Carlo Zaniolo

Computer Science Dept., UCLA, Los Angeles, CA 90095, USA
{ynlaw, zaniolo}@cs.ucla.edu

Abstract. In this paper, we propose an incremental classification algorithm which uses a multi-resolution data representation to find adaptive nearest neighbors of a test point. The algorithm achieves excellent performance by using small classifier ensembles where approximation error bounds are guaranteed for each ensemble size. The very low update cost of our incremental classifier makes it highly suitable for data stream applications. Tests performed on both synthetic and real-life data indicate that our new classifier outperforms existing algorithms for data streams in terms of accuracy and computational costs.

1 Introduction

A significant amount of recent research has focused on mining data streams for applications such as financial data analysis, network monitoring, security, sensor networks, and many others [3,8]. Algorithms for mining data streams have to address challenges not encountered in traditional mining of stored data: at the physical level, these include fast input rates and unending data sets, while, at the logical level, there is the need to cope with concept drift [18]. Therefore, classical classification algorithms must be replaced by, or modified into, incremental algorithms that are fast and light and gracefully adapt to changes in data statistics [17,18,5].

Related Works. Because of their good performance and intuitive appeal, decision tree classifiers and nearest neighborhood classifiers have been widely used in traditional data mining tasks [9]. For data streams, several decision tree classifiers have been proposed—either as single decision trees, or as ensembles of such trees. In particular, VFDT [7] and CVFDT [10] represent well-known algorithms for building single decision tree classifiers, respectively, on stationary, and time-changing data streams. These algorithms employ a criterion based on Hoeffding bounds to decide when a further level of the current decision tree should be created. While this approach assures interesting theoretical properties, the time required for updating the decision tree can be significant, and a large amount of samples is needed to build a classifier with reasonable accuracy. When the size of the training set is small, the performance of this approach can be unsatisfactory.

Another approach to data stream classification uses ensemble methods. These construct a set of classifiers by a base learner, and then combine the predictions

of these base models by voting techniques. Previous research works [17,18,5] have shown that ensembles can often outperform single classifiers and are also suitable for coping with concept drift. On the other hand, ensemble methods suffer from the drawback that they often fail to provide a simple model and understanding of the problem at hand [9].

In this paper, we focus on building nearest neighbor (NN) classifiers for data streams. This technique works well in traditional data mining applications, is supported by a strong intuitive appeal, and it is rather simple to implement. However, the time spent for finding the exact NN can be expensive and, therefore, a significant amount of previous research has focused on this problem. A well-known method for accelerating the nearest neighbor lookup is to use k-d trees [4]. A k-d tree is a balanced binary tree that recursively splits a d-dimensional space into smaller subregions. However, the tree can become seriously unbalanced by massive new arrivals in the data stream, and thus lose the ability of expediting the search. Another approach to NN classifiers attempts to provide approximate answers with error bound guarantees. There are many novel algorithms [11,12,13,14] for finding approximate K -NN on stored data. However, to find the $(1 + \epsilon)$ -approximate nearest neighbors, these algorithms must perform multiple scans of the data. Also, the update cost of the dynamic algorithms [11,13,14] depends on the size of the data set, since the entire data set is needed for the update process. Therefore, they are not suitable for mining data streams.

Our ANNCAD Algorithm. In this paper, we introduce an Addaptive NNearest Classification Algorithm for Data-streams. It is well-known that when data is non-uniform, it is difficult to predetermine K in the KNN classification [6,20]. So, instead of fixing a specific number of neighbors, as in the usual KNN algorithm, we adaptively expand the nearby area of a test point until a satisfactory classification is obtained. To save the computation time for finding adaptive NN, we first preassigning a class to every subregion (cell). To achieve this, we decompose the feature space of a training set and obtain a multi-resolution data representation. There are many decomposition techniques for multi-resolution data representations. The averaging technique used in this paper can be thought of Haar Wavelets Transformation [16]. Thus, information from different resolution levels can then be used for adaptively preassigning a class to every cell. Then we determine to which cell the test point belongs, in order to predict its class. Moreover, because of the compact support property inherited from wavelets, the time spent updating a classifier when a new tuple arrives is a small constant, and it is independent of the size of the data set. Unlike VDFT, which requires a large data set to decide whether to expand the tree by one more level, ANNCAD does not have this restriction.

In the paper, we use grid-based approach for classification. The main characteristic of this approach is the fast processing time and small memory usage, which is independent of the number of data points. It only depends on the number of cells of each dimension in the discretized space, which is easy to adjust in order to fulfill system constraints. Therefore, this approach has been widely employed in clustering problem. Some examples of novel clustering algorithms

are [19], [1] and [15]. However, there is not much work using this approach for classification.

Paper Organization. In this paper, we present our algorithm ANNCAD and discuss its properties in §2. In §3, we compare ANNCAD with some existing algorithms. The results suggest that ANNCAD will outperform existing algorithms. Finally, conclusions and suggestions for future work will be given in §4.

2 ANNCAD

In this section, we introduce our proposed algorithm ANNCAD, which includes four main stages: (1) Quantization of the Feature Space; (2) Building classifiers; (3) Finding predictive label for a test point by adaptively finding its neighboring cells; (4) Updating classifiers for newly arriving tuples. This algorithm only read each data tuple at most once, and only requires a small constant time to process it. We then discuss its properties and complexity.

2.1 Notation

We are given a set of d -dimensional data D with attributes X_1, X_2, \dots, X_d . For each $i = 1, \dots, d$, the domain of X_i is bounded and totally ordered, and ranges over the interval $[L_i, H_i)$. Thus, $X = [L_1, H_1) \times \dots \times [L_d, H_d)$ is the feature space containing our data set D .

Definition 1. Let g be a positive real number, g^d is the volume of a d -dimensional hypercube with side length g . Let $\Delta x_i = (H_i - L_i)/g$ be the i^{th} dimension's quantization interval.

Let B_{i_1, \dots, i_d} denote the block:

$$[L_1 + (i_1 - 1)\Delta x_1, L_1 + i_1\Delta x_1) \times \dots \times [L_d + (i_d - 1)\Delta x_d, L_d + i_d\Delta x_d).$$

Alternatively, we denote B_{i_1, \dots, i_d} by $B_{\mathbf{i}}$, with $\mathbf{i} = (i_1, \dots, i_d)$ the unique identifier for the block. Then, two blocks $B_{\mathbf{k}}$ and $B_{\mathbf{h}}$, $\mathbf{k} \neq \mathbf{h}$, are said to be neighbors if $|k_i - h_i| \leq 1$, for each $i = 1, \dots, d$. In this case, $B_{\mathbf{k}}$ is said to be a neighbor of $B_{\mathbf{h}}$. $Ctr_{B_{\mathbf{i}}}$ denotes the center of block $B_{\mathbf{i}}$, computed as the average of its vertices:

$$Ctr_{B_{\mathbf{i}}} = (L_1 + (i_1 - 1/2)\Delta x_1, \dots, L_d + (i_d - 1/2)\Delta x_d).$$

Definition 2. Let x be a point in the feature space X . Let $B_{\mathbf{i}}$ be a block. The distance between x and $B_{\mathbf{i}}$ is denoted by $dist(x, B_{\mathbf{i}})$.

Note that the distance in Def. 2 can be any kind of distance. In the following, we use Euclidean distance to be the distance between a point and a block.

2.2 Quantization of the Feature Space

The first step of ANNCAD is to partition the feature space into a discretized space with g^d blocks as in Def. 1. It is advisable to choose different sizes of grid according to system resource constraints and desirable fineness of a classifier. For each nonempty block, we count the number of training points contained in it for each class. Now we get the distribution of the data entities in each class. To decide whether we need to start with a finer resolution feature space, we then count the number of training points that do not belong to the majority class of its block as a measure of the training error. We then calculate the coarser representations of the data by averaging the 2^d corresponding blocks in the next finer level. We illustrate the above process by Example 1.

A set of 100 two-class training points in the 2-D unit square is shown in Fig. 1(a). There are two classes for this data set, where a circle (resp. triangle) represents a training point of class I (resp. II). First we separate the training points of each class, discretize them using a 4×4 grid and count the number of training points for each block to get the data distribution of each class (see Fig. 1(b)). Moreover, Fig. 1(c)-(d) show the coarser representations of the data.

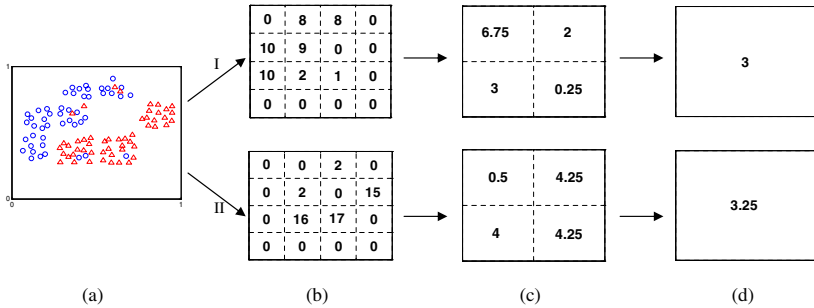


Fig. 1. Multi-resolution representation of a two-class data set

Due to the problem of the curse of dimensionality, the storage amount is exponential in the number of dimensions. To deal with this, we store the nonempty blocks in the leaf nodes of a B^+ -tree using their z-values [21] as keys. Thus the required storage space is much smaller and is bounded by $O(\min(N, g^d))$ where N is the number of training samples. For instance, in Fig. 1, we only need to store information for at most 8 blocks even though there are 100 training points in the 4×4 blocks feature space. To reduce space usage, we may only store the data array of the finest level and calculate the coarser levels on the fly when building a classifier. On the other hand, to reduce time complexity, we may pre-calculate and store the coarser levels. In the following discussion, we assume that the system stores the data representation of each level.

2.3 Building a Classifier and Classifying Test Points

The main idea of ANNCAD is to use a multi-resolution data representation for classification. Notice that the neighborhood relation strongly depends on the quantization process. This will be addressed in next subsection by building several classifier ensembles using different grids obtained by subgrid displacements. Observe that in general, the finer level the block can be classified, the shorter distance between this block and the training set. Therefore, to build a classifier and classify a test point (see Algorithms 1 and 2), we start with the finest resolution for searching nearest neighbors and progressively consider the coarser resolutions, in order to find nearest neighbors adaptively.

We first construct a single classifier as a starting point (see Algorithm 1). We start with setting every block to have a default tag U (Non-visited). In the finest level, we classify any nonempty block with its majority class label. We then classify any nonempty block of every lower level as follows: We label the block by its majority class label if the majority class label has more points than the second majority class by a threshold percentage. If not, we use a specific tag M (Mixed) to label it.

Algorithm 1. BuildClassifier($\{\mathbf{x}, y\}$ | \mathbf{x} is a vector of attributes, y is a class label.)

Quantize the feature space containing $\{\mathbf{x}\}$

Label majority class for each nonempty block in the finest level

For each level $i = \log(g)$ downto 1

For each nonempty block B

If $| \text{majority } c_a | - | \text{2nd majority } c_b | > \text{threshold } \%$, label class c_a

else label tag M

Return Classifier

Algorithm 2. TestClass(test point: \mathbf{t})

For each level $i = \log(g) + 1$ downto 1

If label of $B^i(\mathbf{t}) \ll U$ /* $B^i(\mathbf{t})$ is nonempty */

If label of $B^i(\mathbf{t}) \ll M$, class of \mathbf{t} = class of $B^i(\mathbf{t})$

else class of \mathbf{t} = class of NN of $B^{i+1}(\mathbf{t})$ /* $B^{i+1}(\mathbf{t})$ contains \mathbf{t} in level $i + 1$ */

Break

Return class label for \mathbf{t} , $B^i(\mathbf{t})$

We build a classifier for the data set of Example 1 and set the threshold value to be 80%. Fig. 2(a), (b) and (c) show the class label of each nonempty block in the finest, intermediate and coarsest resolution respectively.

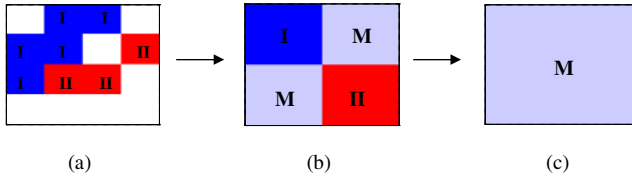


Fig. 2. Hierarchical structure of classifiers

For each level i , a test point \mathbf{t} belongs to a unique block $B^i(\mathbf{t})$. We search from the finest to the coarsest level until reaching a nonempty block $B^i(\mathbf{t})$. If the label of $B^i(\mathbf{t})$ is one of the classes, we label the test point by this class. Otherwise, if $B^i(\mathbf{t})$ has tag M , we find the nearest neighbor block of $B^{i+1}(\mathbf{t})$ where $B^{i+1}(\mathbf{t})$ is a block containing \mathbf{t} in level $i + 1$. To reduce the time spent, we only consider the neighbors of $B^{i+1}(\mathbf{t})$ which belong to $B^i(\mathbf{t})$ in level i . It is very easy to access these neighbors as they are also neighbors of $B^{i+1}(\mathbf{t})$ in the B^+ -tree with their z -values as keys. Note that $B^{i+1}(\mathbf{t})$ must be empty, otherwise we should classify it at level $i + 1$. But some of the neighbors of $B^{i+1}(\mathbf{t})$ must be nonempty as $B^i(\mathbf{t})$ is nonempty. We simply calculate the distance between test point \mathbf{t} and each neighbor of $B^{i+1}(\mathbf{t})$ and label \mathbf{t} by the class of NN.

We use the classifier built in Example 2 to classify a test point $\mathbf{t} = (0.6, 0.7)$. Starting with the finest level, we found that the first nonempty block containing \mathbf{t} is $[0.5, 1) \times [0.5, 1)$ (see Fig. 3(b)). Since it has tag M , we calculate the distance between \mathbf{t} and each nonempty neighboring block in the next finer level ($[0.75, 1) \times [0.5, 0.75), [0.5, 0.75) \times [0.75, 1)$). Finally, we get the nearest neighboring block $[0.75, 1) \times [0.5, 0.75)$ and label \mathbf{t} to be class I (see Fig.

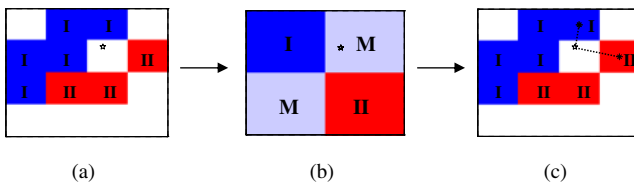


Fig. 3. Hierarchical classifier access

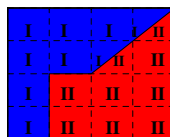


Fig. 4. The combined classifier

3(c)). When we combine the multi-resolution classifier of each level, we get a classifier for the whole feature space (see Fig. 4).

2.4 Incremental Updates of Classifiers

The main requirement of a data stream classification algorithm is that it is able to update classifiers incrementally and effectively when a new tuple arrives. Moreover, updated classifier should be adapt to concept drift behaviors. In this subsection, we present incremental update process of ANNCAD for a stationary data, without re-scanning the data and discuss an exponential forgetting technique to adapt to concept drifts.

Because of the compact support property, arrival of a new tuple only affects the blocks of the classifier in each level containing this tuple. Therefore, we only need to update the data array of these blocks and their classes if necessary. During the update process, the system may run out of memory as the number of nonempty blocks may increase. To deal with this, we may simply remove the finest data array, multiple the entries of the remaining coarser data arrays by 2^d , and update the quantity g . A detailed description of updating classifiers can be found in Algorithm 3. This solution can effectively meet the memory constraint.

Algorithm 3. UpdateClassifier(new tuple: \mathbf{t})

For each level $i = \log(g) + 1$ **downto** 1
 Add $\delta_t / 2^{d \times (\log(g) + 1 - i)}$ to data array Φ^i
 /* δ_t is a matrix with value 1 in the corr. entry of t and 0 elsewhere.*/
 If i is the finest level, label $B^i(\mathbf{t})$ with the majority class
 else if $|\text{majority } c_a| - |\text{2nd majority } c_b| > \text{threshold } \%$, label $B^i(\mathbf{t})$ by c_a
 else label $B^i(\mathbf{t})$ by tag M
If memory runs out,
 Remove the data array of level $\log(g) + 1$
 For each level $i = \log(g)$ **downto** 1, $\Phi^i = 2^d \cdot \Phi^i$
 Label each nonempty block of the classifier in level $\log(g)$ by its majority class
 Set $g = g/2$
Return updated classifier

Exponential Forgetting. If the concept of the data changes over time, a very common technique called exponential forgetting may be used to assign less weight to the old data to adapt to more recent trend. To achieve this, we multiply an exponential forgetting factor λ to the data array, where $0 \leq \lambda \leq 1$. For each level i , after each time interval t , we update the data array Φ^i to be:

$$\Phi^i|_{(n+1)t} \leftarrow \lambda \Phi^i|_{n \cdot t}$$

where $\Phi^i|_{n \cdot t}$ is the data array at time $n \cdot t$. Indeed, if there is no concept change, the result of classifier will not be affected. If there is a concept drift, the classifier

can adapt to the change quickly since the weight of the old data is exponentially decreased. In practice, an exponential forgetting technique is easier to implement than a sliding window because we need extra memory buffer to store the data of the most current window for implementing the sliding window.

2.5 Building Several Classifiers Using Different Grids

As mentioned above, the neighborhood relation strongly depends on the quantization process. For instance, consider the case that there is a training point u which is close to the test point v but they are located in different blocks. Then the information on u may not affect the classification of v .

To overcome the problem of initial quantization process, we build several classifier ensembles starting with different quantization space. In general, to build n^d different classifiers, each time we shift $\frac{1}{n}$ of the unit length of feature space for a set of selected dimensions. Fig. 5 shows a reference grid and its 3 different shifted grids for a feature space with 4×4 blocks. For a given test point \mathbf{t} , we use these n^d classifiers to get n^d class labels and selected blocks $B^i(\mathbf{t})$ of \mathbf{t} in each level i , starting from the finest one. We then choose the majority class label. If there is tie, we calculate the distance between each selected block $B^i(\mathbf{t})$ with majority class label and \mathbf{t} to find the closest one. Algorithm 4 shows this classifying process using n^d classifiers.

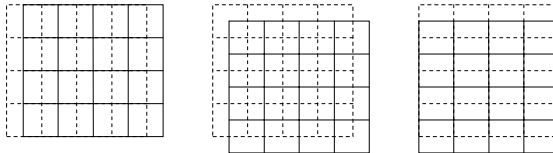


Fig. 5. An example of 4 different grids for building 4 classifiers

Algorithm 4. Test n^d Class(objects: \mathbf{t})

For each level $i = \log(g) + 1$ downto 1

Get the label of \mathbf{t} for each classifier

If there is a label $\langle \rangle U$, choose the majority label

If there is a tie, label \mathbf{t} by class of $B^i(\mathbf{t})$ with closest center to \mathbf{t}

Break

Return class label for \mathbf{t}

The following theorem shows that the approximation error of finding nearest neighbors decreases as the number of classifier ensembles increases.

Theorem 1. Let $x \in X$ and Y be a set of n^d classifiers. If $z \notin Y$ and $(x, y) < (1 + \frac{1}{n-1})^* \dots (x, z)$ for $y \in Y$

For simplicity, we consider the case when $d = 1$. This proof works for any d . For $d = 1$, we build n classifiers, where each classifier i use the grid that is shifted $\frac{i}{n}$ unit length from the original grid. Let ϵ be the length of a block. Consider a test point x , x belongs to an interval I_k for classifier k . Note that $[x - \frac{n-1}{n}\epsilon, x + \frac{n-1}{n}\epsilon] \subset \bigcup I_k \subset [x - \epsilon, x + \epsilon]$. Hence, the distance between x and its nearest neighbor that we found must be less than ϵ . Meanwhile, the points that we do not consider should be at least $\frac{n-1}{n}\epsilon$ far away from x . If $z \notin Y$, $\frac{dist(x,y)}{dist(x,z)} < \frac{\epsilon}{(n-1)\epsilon/n} = (1 + \frac{1}{n-1})$ for every $y \in Y$.

The above theorem shows that the classification result using one classifier does not have any guarantee about the quality of the nearest neighbors that it found because the ratio of approximation error will tend to infinity. When n is large enough, the set of training points selected by those classifier ensembles are exactly the set of training points with distance ϵ from the test point. To achieve an approximation error bound guarantee, theoretically we need an exponential number of classifiers. However, in practice, we only use two classifiers to get a good result. Indeed, experiments in §3 show that few classifiers can obtain a significant improvement at the beginning. After this stage, the performance will become steady even though we keep increasing the number of classifiers.

2.6 Properties of ANNCAD

As ANNCAD is a combination of multi-resolution and adaptive nearest neighbors techniques, it inherits both their properties and their advantages.

- **Locality**. The locality property allows a fast update. As a new tuple arrival only affects the class of the block containing it in each level, the incremental update process only costs a constant time (number of levels).
- **Threshold**. We may set a threshold value for classifying decisions to remove noise.
- **Multi-resolution**. This algorithm makes it easy to build multi-resolution classifiers. Users can specify the number of levels to efficiently control the fineness of the classifier. Moreover, one may optimize the system resource constraints and easy to adjust on the fly when the system runs out of memory.
- **Complexity**. Let g , N and d be the number of blocks of each dimension, training points and attributes respectively. The time spent on building a classifier is $O(\min(N, g^d))$ with constant factor $\log(g)$. For the time spent on classifying a test point, the worst case complexity is $O(\log_2(g) + 2^d)$ where the first part is for classifying a test point using classifiers and the second part is for finding its nearest neighbor which is optional. Also, the time spent for updating classifiers when a new tuple arrives is $\log_2(g) + 1$. Comparing with the time spent in VFDT, our method is more attractive.

3 Performance Evaluation

In this section, we first study the effects on parameters for ANNCAD by using two synthetic data sets. We then compare ANNCAD with VFDT and CVFDT

on three real-life data sets. To illustrate the approximation power of ANNCAD, we include the results of \dots , which computes ANN exactly, as controls. \dots : For each test point t , we search the area within 0.5 block side length distance. If the area is nonempty, we classify t as the majority label of all these points in this area. Otherwise, we expand the searching area by doubling the radius until we get a class for t . Note that the time and space complexities of \dots are very expensive making it impractical to use.

3.1 Synthetic Data Sets

The aim of this experiment is to study the effect on the initial resolution for ANNCAD. In this synthetic data set, we consider a 3-D unit cube. We randomly pick 3k training points and assign those points which are inside a sphere with center (0.5, 0.5, 0.5) and radius 0.5 to be class 0, and class 1 otherwise. This data set is effective to test the performance of a classifier as it has a curve-like decision boundary. We then randomly draw 1k test points and run ANNCAD starting with different initial resolution and 100% threshold value. In Fig. 6(a), the result shows that a finer initial resolution gets a better result. This can be explained by the fact that we can capture a curve-like decision boundary if we start with a finer resolution. On the other hand, as discussed in last section, the time spent for building a classifier increases linearly for different resolutions. In general, we should choose a resolution according to system resource constraints.

The aim of this experiment is to study the effect on number of classifier ensembles for ANNCAD. As in the previous experiment, we randomly pick 1k training examples and assign them labels. We then randomly draw 1k test points and test them based on the voting result of these classifiers. We set $16 \times 16 \times 16$ blocks for the finest level and 100% threshold value. In Fig. 6(b), the result shows that having more classifiers will get a better result in the beginning. The performance improvement becomes steady even though we keep increasing the number of classifiers. It is because there is no further information given when we increase the number of classifiers. In this experiment, we only use 2 or 3 classifiers to obtain a competitive result with the \dots (90.4%).

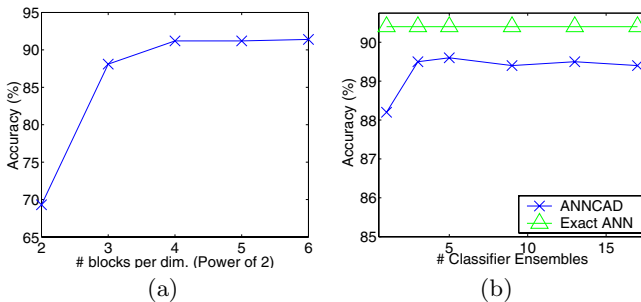


Fig. 6. Effect on initial resolutions and number of classifiers

3.2 Real Life Data Sets

The aim of this set of experiments is to compare the performance of ANNCAD with that of VFDT and CFVDT on stationary and time-changing real-life data sets respectively. We first used a letter recognition data set from the UCI machine learning repository web site [2]. The objective is to identify a black-and-white pixel displays as one of the 26 English alphabet. In this data set, each entity is a pixel display for an English alphabet and has 16 numerical attributes to describe its pixel displays. The detail description of this data set is provided in [2]. In this experiment, we use 15k tuples for training set with 5% noise added and 5k for test set. We obtain noisy data by randomly assigning a class label for 5% training examples. For ANNCAD, we set g for the initial grid to be 16 units and build two classifiers. Moreover, since VFDT needs a very large training set to get a fair result, we rescan the data sets up to 500 times for VFDT. So the data set becomes 7,500,000 tuples. In Fig. 7(a), the performance of ANNCAD dominates that of VFDT. Moreover, ANNCAD only needs one scan to achieve this result, which shows that ANNCAD even works well for a small training set.

The second real life data set we used is the Forest Cover Type data set which is another data set from [2]. The objective is to predict forest cover type (7 types). For each observation, there are 54 variables. Neural network (backpropagation) was employed to classify this data set and got 70% accuracy, which is the highest one recorded in [2]. In our experiment, we used all the 10 quantitative variables. There are 12k examples for training set and 90k examples for testing set. For ANNCAD, we scaled each attribute to the range $[0, 1)$. We set g for the initial grid to be 32 units and build two classifiers. As the above experiment, we rescan the training set up to 120 times for VFDT, until its performance becomes steady. In Fig. 7(b), the performance of ANNCAD dominates that of VFDT. These two experiments show that ANNCAD works well in different kinds of data sets.

We further tested ANNCAD in the case when there are concept drifts in data set. The data we used was extracted from the census bureau database [2]. Each observation represents a record of an adult and has 14 attributes including age, race etc. The prediction task is to determine whether a person makes over 50K a year. Concept drift is simulated by grouping records with same race (Amer-Indian-Eskimo(AIE), Asian-Pac-Islander(API), Black(B), Other(O), White(W)). The distribution of training tuples of each race is shown in Fig. 7(c). Since the models for different races of people should be different, concept drifts are introduced when $n = 311, 1350, 4474, 4746$. In this experiment, we used the 6 continuous attributes. We used 7800 examples for learning and tested the classifiers for every 300 examples. For ANNCAD, we build two classifiers and set λ to be 0.98 and g for the initial grid to be 64 units. We scaled the attribute values as mentioned in the previous experiment. The results are shown in Fig. 7(c). The curves show that ANNCAD keeps improving in each region. Also, as mentioned in §2.6, computations required for ANNCAD are much lower than CVFDT.

Moreover, notice that ANNCAD works almost as well as on these three data sets, which demonstrates its excellent approximation ability.

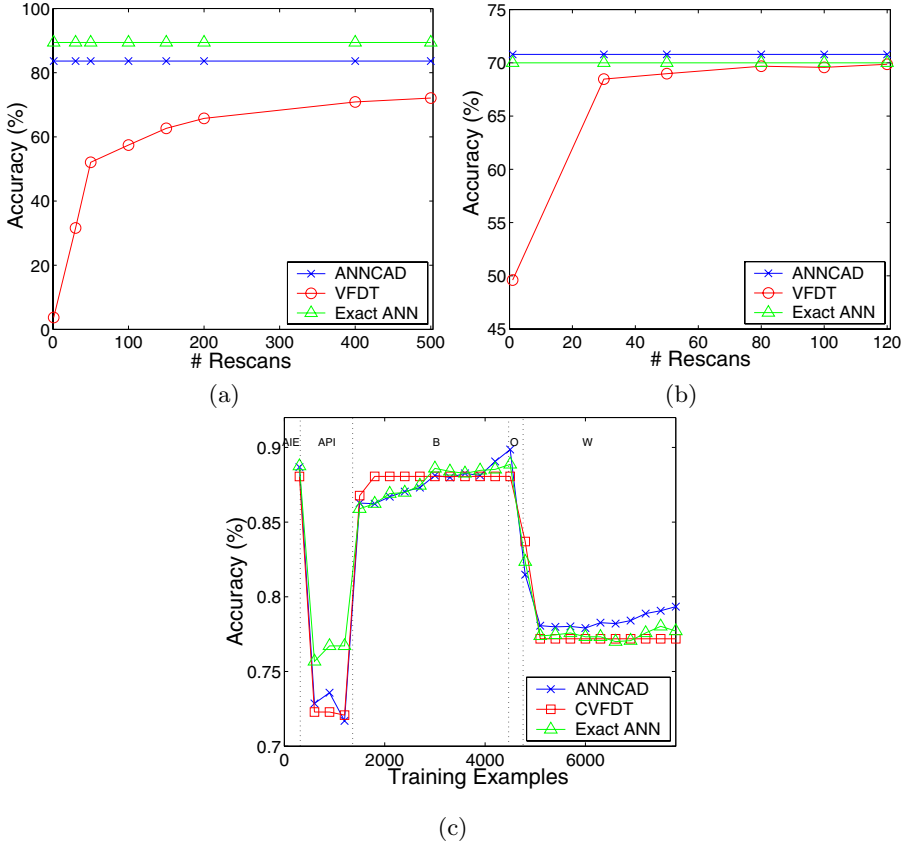


Fig. 7. Three real-life data sets:(a) Letter Recognition (b) Forest Covertype (c) Census

4 Conclusion and Future Work

In this paper, we proposed an incremental classification algorithm ANNCAD using a multi-resolution data representation to find adaptive nearest neighbors of a test point. ANNCAD is very suitable for mining data streams as its update speed is very fast. Also, the accuracy compares favorably with existing algorithms for mining data streams. ANNCAD adapts to concept drift effectively by the exponential forgetting approach. However, the very detection of sudden concept drift is of interest in many applications. The ANNCAD framework can also be extended to detect concept drift—e.g. changes in class label of blocks is a good indicator of possible concept drift. This represents a topic for our future research.

Acknowledgement. This research was supported in part by NSF Grant No. 0326214.

References

1. R. Agrawal, J. Gehrke, D. Gunopulos and P. Raghavan. Automatic Subspace Clustering of High Dimensional Data for Data Mining Applications. *SIGMOD 1998*: 94–105.
2. C. L. Blake and C. J. Merz. UCI Repository of machine learning databases. <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
3. B. Babcock, S. Babu, R. Motawani and J. Widom. Models and issues in data stream systems. *PODS 2002*: 1–16.
4. J. Bentley. Multidimensional binary search trees used for associative searching. *Communication of the ACM 18(9)*: 509–517 (1975).
5. F. Chu and C. Zaniolo. Fast and light boosting for adaptive mining of data streams. *PAKDD 2004*: 282–292.
6. C. Domeniconi, J. Peng and D. Gunopulos, Locally adaptive metric nearest-neighbor classification, *IEEE Transactions on Pattern Analysis and Machine Intelligence 24(9)*: 1281–1285 (2002).
7. P. Domingos and G. Hulten. Mining high-speed data streams. *KDD 2000*: 71–80.
8. L. Golab and M. Özsu. Issues in data stream management. *ACM SIGMOD 32(2)*: 5–14 (2003).
9. J. Han and M. Kamber. *Data Mining – Concepts and Techniques (2000)*. Morgan Kaufmann Publishers.
10. G. Hulten, L. Spence and P. Domingos. Mining time-changing data streams. *KDD 2001*: 97–106.
11. P. Indyk, R. Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. *STOC 1998*: 604–613.
12. P. Indyk. Dimensionality reduction techniques for proximity problems. *ACM-SIAM symposium on Discrete algorithms 2000*: 371–378.
13. P. Indyk. High-dimensional computational geometry. *Dept. of Comput. Sci., Stanford Univ., 2001*.
14. E. Kushilevitz, R. Ostrovsky, Y. Rabani. Efficient Search for Approximate Nearest Neighbor in High Dimensional Spaces. *SIAM J. Comput. 30(2)*: 457–474 (2000).
15. G. Sheikholeslami, S. Chatterjee and A. Zhang. WaveCluster: A Multi-Resolution Clustering Approach for Very Large Spatial Databases. *VLDB 1998*: 428–439.
16. G. Strang and T. Nguyen. *Wavelets and Filter Banks (1996)*. Wellesley-Cambridge Press.
17. W. Street and Y. Kim. A streaming ensemble algorithm (SEA) for large-scale classification. *SIGKDD 2001*: 377–382.
18. H. Wang, W. Fan, P. Yu and J. Han. Mining concept-drifting data streams using ensemble classifiers. *SIGKDD 2003*: 226–235.
19. W. Wang, J. Yang and R. Muntz. STING: A Statistical Information Grid Approach to Spatial Data Mining *VLDB 1997*: 186–195.
20. D. Wettschereck and T. Dietterich. Locally Adaptive Nearest Neighbor Algorithms. *Advances in Neural Information Processing Systems 6*: 184–191 (1994).
21. C. Zaniolo, S. Ceri, C. Faloutsos, R. Snodgrass, V. Subrahmanian and R. Zicari. *Advanced Database Systems (1997)*. Morgan Kaufmann Press.

Support Vector Random Fields for Spatial Classification

Chi-Hoon Lee, Russell Greiner, and Mark Schmidt

Department of Computing Science,
University of Alberta,
Edmonton AB, Canada
{chihoon, greiner, schmidt}@cs.ualberta.ca

Abstract. In this paper we propose Support Vector Random Fields (SVRFs), an extension of Support Vector Machines (SVMs) that explicitly models spatial correlations in multi-dimensional data. SVRFs are derived as Conditional Random Fields that take advantage of the generalization properties of SVMs. We also propose improvements to computing posterior probability distributions from SVMs, and present a local-consistency potential measure that encourages spatial continuity. SVRFs can be efficiently trained, converge quickly during inference, and can be trivially augmented with kernel functions. SVRFs are more robust to class imbalance than Discriminative Random Fields (DRFs), and are more accurate near edges. Our results on synthetic data and a real-world tumor detection task show the superiority of SVRFs over both SVMs and DRFs.

1 Introduction

The task of classification has traditionally focused on data that is “independent and identically distributed” (iid), in particular assuming that the class labels for different data points are conditionally independent (ie. knowing that one patient has cancer does not mean another one will). However, real-world classification problems often deal with data points whose labels are correlated, and thus the data violates the iid assumption. There is extensive literature focusing on the 1-dimensional ‘sequential’ case (see [1]), where correlations in the labels of data points in a linear sequence exist, such as in strings, sequences, and language. This paper focuses on the more general ‘spatial’ case, where these correlations exist in data with two-dimensional (or higher-dimensional) structure, such as in images, volumes, graphs, and video.

Classifiers that make the iid assumption often produce undesirable results when applied to data with spatial dependencies in the labels. For example, in the task of image labeling, a classifier could classify a pixel as ‘face’, even if all adjacent pixels were classified as ‘non-face’. This problem motivates the use of Markov Random Fields (MRFs) and more recently Conditional Random Fields (CRFs) for spatial data. These classification techniques augment the performance of an iid classification technique (often a Mixture Model for MRFs, and Logistic Regression for CRFs) by taking into account spatial class dependencies.

Support Vector Machines (SVMs) are classifiers that have appealing theoretical properties [2], and have shown impressive empirical results in a wide variety of tasks. However, this technique makes the critical iid assumption. This paper proposed an extension to SVMs that considers spatial correlations among data instances (as in Random Field models), while still taking advantage of the powerful discriminative properties of SVMs. We refer to this technique as Support Vector Random Fields (SVRFs)

The remaining sections of this paper are organized as follows. Section 2 formalizes the task and reviews related methods for modeling dependencies in the labels of spatial data. Section 3 reviews Support Vector Machines, and presents our Support Vector Random Field extension. Experimental results on synthetic and real data sets are given in Sect. 4, while a summary of our contribution is presented in Sect. 5.

2 Related Work

The challenge of performing classification while modeling class dependencies is often divided into two perspectives: Generative and Discriminative models [1]. Generative classifiers learn a model of the joint probability, $p(x, y) = p(x|y)p(y)$, of the features x and corresponding labels y . Predictions are made using Bayes rule to compute $p(y|x)$, and finding an assignment of labels maximizing this probability. In contrast, discriminative classifiers model the posterior $p(y|x)$ directly without generating any prior distributions over the classes. Thus, discriminative models solely focus on maximizing the conditional probability of the labels, given the features. For many applications, discriminative classifiers often achieve higher accuracy than generative classifiers [1]. There has been much related work on using random field theory to model class dependencies in generative and more recently discriminative contexts [3,4]. Hence, we will first review [\[3\]](#) (typically formulated as a generative classifier), followed by [\[4\]](#) (a state-of-the-art discriminative classifier built upon the foundations of Markov Random Fields).

2.1 Problem Formulation

In this work, we will focus on the task of classifying elements (pixels or regions) of a two-dimensional image, although the methods discussed also apply to higher-dimensional data. An image is represented with an M by N matrix of elements. For an instance $X = (x_{11}, x_{12}, \dots, x_{1N}, \dots, x_{M1}, x_{M2}, \dots, x_{MN})$, we seek to infer the most likely joint class labels:

$$Y^* = (y_{11}^*, y_{12}^*, \dots, y_{1N}^*, \dots, y_{M1}^*, y_{M2}^*, \dots, y_{MN}^*)$$

If we assume that the labels assigned to elements are independent, the following joint probability can be formulated: $P(Y) = \prod_{i=1}^M \prod_{j=1}^N P(y_{ij})$. However, conditional independency does not hold for image data, since spatially adjacent elements are likely to receive the same labels. We therefore need to explicitly

consider this local dependency. This involves addressing three important issues: How should the optimal solution be defined, how are spatial dependencies considered, and how should we search the (exponential size) configuration space.

2.2 Markov Random Fields (MRFs)

Markov Random Fields (MRFs) provide a mathematical formulation for modeling local dependencies, and are defined as follows [3]:

Definition 1.

A Markov Random Field is a triplet (S, N, P) where $S = \{1, \dots, n\}$ is a finite set of nodes, $N = \{N_i \mid i \in S\}$ is a set of neighborhoods, and $P = \{P(y_i | y_{S-\{i\}}) \mid i \in S\}$ is a set of conditional distributions.

$$P(Y) > 0$$

$$P(y_i | y_{S-\{i\}}) = P(y_i | y_N)$$

Condition 2 (Markovianity) states that the conditional distribution of an element y_i is dependent only on its neighbors. Markov Random Fields have traditionally sought to maximize the joint probability $P(Y^*)$ (a generative approach). In this formulation, the posterior over the labels given the observations is formulated using Bayes' rule as:

$$P(Y|X) \propto P(X|Y)P(Y) = P(Y) \prod_i^n P(x_i|y_i) \quad (1)$$

In (1), the equivalence between MRFs and Gibbs Distributions [5] provides an efficient way to factor the prior $P(Y)$ over cliques defined in the neighborhood Graph G . The prior $P(Y)$ is written as

$$P(Y) = \frac{\exp(\sum_{c \in C} V_c(Y))}{\sum_{Y' \in \Omega} \exp(\sum_{c \in C} V_c(Y'))} \quad (2)$$

where $V_c(Y)$ is a clique potential function of labels for clique $c \in C$, C is a set of cliques in G , and Ω is the space of all possible labelings. From (1) and (2), the target configuration Y^* is a realization of a locally dependent Markov Random Field with a specified prior distribution. Based on (1) and (2) and using Z to denote the (normalizing) "partition function", if we assume Gaussian likelihoods then the posterior distribution can be factored as:

$$P(Y|X) = \frac{1}{Z} \exp \left[\sum_{i \in S} \log(P(x_i|y_i)) + \sum_{c \in C} V_c(Y_c) \right] \quad (3)$$

The Gaussian assumption for $P(X|Y)$ in (1) allows straightforward Maximum Likelihood parameter estimation. Although there have been many approximation

algorithms designed to find the optimal Y^* , we will focus on a local method called [\[5\]](#), written as:

$$y_i^* = \arg \max_{y \in L} P(y_i | y_N, x_i) \tag{4}$$

Assuming Gaussians for the likelihood and a pairwise neighborhood system for the prior over labels, (4) can be restated as:

$$y_i^* = \arg \max_{y \in L} \frac{1}{Z_i} \exp \left[\log(P(x_i | y_i)) + \sum_{j \in N} \beta y_i y_j \right] \tag{5}$$

where β is a constant and L is a set of class labels.

This concept has proved to be applicable in a wide variety of domains where there exists correlations among neighboring instances. However, the generative nature of the model and the assumption that the likelihood is Gaussian can be too restrictive to capture complex dependencies between neighboring elements or between observations and labels. In addition, the prior over labels is completely independent from the observations, thus the interactions between neighbors are not proportional to their similarity.

2.3 Conditional Random Fields (CRFs)

to avoid the Gaussian assumption by using a model that seeks to maximize the conditional probability of the labels given the observations $P(Y^* | X)$ (a discriminative model), and are defined as follows [\[1\]](#):

Definition 2. Let $G = (S, E)$ be an undirected graph with nodes S and edges E . Let Y be a set of labels for the nodes S . Let X be a set of observations for the nodes S . Let y_i be the label of node i . Let x_i be the observation of node i . Let $y_{S \setminus i}$ be the labels of all nodes except i . Let $x_{S \setminus i}$ be the observations of all nodes except i . Let $P(y_i | X, y_{S \setminus i}) = P(y_i | X, y_N)$ be the conditional probability of the label y_i given the observations X and the labels of all other nodes y_N .

This model alleviates the need to model the observations $P(X)$, allowing the use of arbitrary attributes of the observations without explicitly modeling them. CRFs assume a 1-dimensional chain-structure where only adjacent elements are neighbors. This allows the factorization of the joint probability over labels. Discriminative Random Fields (DRFs) extend 1-dimensional CRFs to 2-dimensional structures [\[6\]](#). The conditional probability of the labels Y in the Discriminative Random Field framework is defined as:

$$P(Y | X) = \frac{1}{Z} \exp \left(\sum_{i \in S} A_i(y_i, X) + \sum_{i \in S} \sum_{j \in N} I_{ij}(y_i, y_j, X) \right) \tag{6}$$

A_i is the ‘Association’ potential that models dependencies between the observations and the class labels, while I_i is the ‘Interaction’ potential that models dependencies between the labels of neighboring elements (and the observations). Note that this is a much more powerful model than the assumed Gaussian Association potential and the indicator function used for the Interaction potential

(that doesn't consider the observations) in MRFs. Parameter learning in DRFs involves maximizing the log likelihood of (6), while inference uses ICM [6].

DRFs are a powerful method for modeling dependencies in spatial data. However, several problems associated with this method include the fact that it is hard to find a good initial labeling and stopping criteria during inference, and it is sensitive to issues of class imbalance. Furthermore, for some real-world tasks the use of logistic regression as a discriminative method in DRFs often does not produce results that are as accurate as powerful classification models such as Support Vector Machines (that make the iid assumption).

3 Support Vector Random Fields (SVRFs)

This section presents Support Vector Random Fields (SVRFs), our extension of SVMs that allows the modelling of non-trivial two-dimensional (or higher) spatial dependencies using a CRF framework. This model has two major components: The *observation-matching potential function* and the *local-consistency potential function*. The *observation-matching potential function* captures relationships between the observations and the class labels, while the *local-consistency potential function* models relationships between the labels of neighboring data points and the observations at data points. Since the selection of the observation-matching potential is critical to the performance of the model, the Support Vector Random Field model employs SVMs for this potential, providing a theoretical and empirical advantage over the logistic model used in DRFs and the Gaussian model used in MRFs, that produce unsatisfactory results for many tasks. SVRFs can be formulated as follows:

$$P(Y|X) = \frac{1}{Z} \exp \left\{ \sum_{i \in S} \log(O(y_i, \mathcal{Y}_i(X))) + \sum_{i \in S} \sum_{j \in N} V(y_i, y_j, X) \right\} \quad (7)$$

In this formulation, $\mathcal{Y}_i(X)$ is a function that computes features from the observations X for location i , $O(y_i, \mathcal{Y}_i(X))$ is the observation-potential, and $V(y_i, y_j, X)$ is the local-consistency potential. The pair-wise neighborhood system is defined as a local dependency structure. In this work, interactions between pixels with a Euclidean distance of 1 were considered (ie. the radius 1 von Neumann neighborhood). We will now examine these potentials in more detail.

3.1 Observation-Matching

The observation-matching potential seeks to find a posterior probability distribution that maps from the observations to corresponding class labels. DRFs employ a Generalized Linear Models (GLM) for this potential. However, GLMs often do not estimate appropriate parameters. This is especially true in image data where feature sets may have a high number of dimensions and/or several features have a high degree of correlation. This can cause problems in parameter estimation and approximations to resolve these issues may not produce optimal parameters [7].

Fortunately, the CRF framework allows a flexible choice of the observation-matching potential function. We overcome the disadvantages of the GLM by employing a Support Vector Machine classifier, seeking to find the margin maximizing hyperplane between the classes. This classifier has appealing properties in high-dimensional spaces and is less sensitive to class imbalance. Furthermore, due to the properties of error bounds, SVMs tends to outperform GLMs, especially when the classes overlap in the feature space (often the case with image data). Parameter estimation for SVMs involves optimizing the following Quadratic Programming problem for the training data x_i (where, C is a constant that bounds the misclassification error):

$$\begin{aligned} \max \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_i^N \sum_j^N \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{subject to } 0 \leq \alpha_i \leq C \text{ and } \sum_{i=1}^N \alpha_i y_i = 0 \end{aligned} \quad (8)$$

Consequently, the decision function, given the parameters α_i for the l training instances and bias term b , is (for a more thorough discussion of SVMs, we refer to [2]): $f(x) = \sum_{i=1}^l (\alpha_i y_i x \cdot x_i) + b$

Unfortunately, the decision function $f(x)$ produced by SVMs measures distances to the decision boundary, while we require a posterior probability function. We adopted the approach of [8] to convert the decision function to a posterior probability function. This approach is efficient and minimizes the risk of overfitting during the conversion, but has some ambiguities and potential difficulties in numerical computation. We have addressed these issues in our approach, which will be briefly outlined here.

We estimate a posterior probability from the Support Vector Machine decision function using the sigmoid function:

$$O(y_i = 1, \mathcal{Y}_i(X)) = \frac{1}{1 + \exp(Af(\mathcal{Y}_i(X)) + B)} \quad (9)$$

The parameters A and B are estimated from training data represented as pairs $(f(\mathcal{Y}_i(X)), t_i)$, where $f(\cdot)$ is the Support Vector Machine decision function, and t_i denotes a relaxed probability that $y_i = 1$ as in (9). We could set $t_i = 1$, if the class label at i is 1 (ie. $y_i = 1$). However, this ignores the possibility that $\mathcal{Y}_i(X)$ has the opposite class label (ie. -1). Thus, we employed the relaxed probability: $t_i = \frac{N_+ + 1}{N_+ + 2}$, if $y_i = 1$, and $t_i = \frac{1}{N_- + 2}$, if $y_i = -1$ (N_+ and N_- being the number of positive and negative class instances). By producing the new forms of training instances, we can solve the following optimization problem to estimate parameters:

$$\min - \sum_{i=1}^l \left[t_i \log p(\mathcal{Y}_i(X)) + (1 - t_i) \log(1 - p(\mathcal{Y}_i(X))) \right] \quad (10)$$

where

$$p(\mathcal{Y}_i(X)) = \frac{1}{1 + \exp(Af(\mathcal{Y}_i(X)) + B)}$$

[8] adopted a Levenberg-Marquardt approach to solve the optimization problem, finding an approximation of the Hessian matrix. However, this may cause incorrect computations of the Hessian matrix (especially for unconstrained optimizations [7]). Hence, we employed Newton's method with backtracking line search to solve the optimization. In addition, in order to avoid overflows and underflows of *exp* and *log* functions, we reformulate Eq.10 as follows:

$$\begin{aligned} & -\left(t_i \log p(\mathcal{Y}_i(X)) + (1 - t_i) \log(1 - p(\mathcal{Y}_i(X)))\right) \\ & = t_i(Af(\mathcal{Y}_i(X)) + B) + \log(1 + \exp(-Af(\mathcal{Y}_i(X)) - B)) \end{aligned} \quad (11)$$

3.2 Local-Consistency

In MRFs, local-consistency considers correlations between neighboring data points, and is considered to be observation independent. CRFs provide more powerful modelling of local-consistency by removing the assumption of observation independence. In order to use the principles of CRFs for local-consistency, an approach is needed that penalizes discontinuity between pairwise sites. For this, we use a linear function of pairwise continuity:

$$V(y_i, y_j, X) = y_i y_j \nu^T \Phi_{ij}(X) \quad (12)$$

$\Phi_{ij}(X)$ is a function that computes features for sites i and j based on observations X . As opposed to DRFs, which penalize discontinuity by considering the absolute difference between pairwise observations [6], our approach introduces a new mapping function $\Phi(\cdot)$ that encourages continuity in addition to penalizing discontinuity (using $\max(\mathcal{Y}(X))$ to denote the vector of max values for each feature):

$$\Phi_{ij}(X) = \frac{\max(\mathcal{Y}(X)) - |\mathcal{Y}_i(X) - \mathcal{Y}_j(X)|}{\max(\mathcal{Y}(X))} \quad (13)$$

3.3 Learning and Inference

The proposed model needs to estimate the parameters of the observation-matching function and the local-consistency function. Although we estimate these parameters sequentially, our model outperforms the simultaneous learning approach of DRFs and significantly increases its computational efficiency.

The parameters of the Support Vector Machine decision function are first estimated by solving the Quadratic Programming problem in (8) (using SVMlight [9]). We then convert the decision function to a posterior function using (10) and the new training instances. Finally, we adopted pseudolikelihood [3] to estimate the local consistency parameters ν , due to its simplicity and fast computation. For training on l pixels from K images, pseudolikelihood is formulated as:

$$\hat{\nu} = \arg \max_{\nu} \prod_{k=1}^K \prod_{i=1}^l P(y_i^k | y_N^k, X^k, \nu) \quad (14)$$

As in [6], to ensure that the log-likelihood is convex we assume that ν is Gaussian and compute the local-consistency parameters using its log likelihood $l(\hat{\nu})$:

$$l(\hat{\nu}) = \arg \max_{\nu} \sum_{k=1}^K \sum_{i=1}^l \left\{ O_i^n + \sum_{j \in N} V(y_i^k, y_j^k, X^k) - \log(z_i^k) \right\} - \frac{1}{2\tau} \nu^T \nu \quad (15)$$

In this model, z_i^k is a partition function for each site i in image k , and τ is a regularizing constant. Equation (15) is solved by gradient descent, and note that the observation matching function acts as a constant during this process. Due to the employment of SVMs, the time complexity of learning is $O(S^2)$, where S is the number of pixels to be trained, although in practice it is much faster.

The inference problem is to infer an optimal labeling Y^* given a new instance X and the estimated model parameters. We herein adopted the Iterated Conditional Modes (ICM) approach described in Section 2.2 [5], that maximizes the local conditional probability iteratively. For our proposed model and [6], ICM is expressed as,

$$y_i^* = \arg \max_{y \in L} P(y_i | y_N, X) \quad (16)$$

Although ICM is based on iterative principles, it often converges quickly to a high quality configuration, and each iteration has time complexity $O(S)$.

4 Experiments

We have evaluated our proposed model on synthetic and real-world binary image labeling tasks, comparing our approach to Logistic Regression, SVMs, and DRFs for these problems. Since class imbalance was present in many of the data sets, we used the Jaccard measure to quantify performance: $f = \frac{TP}{TP+FP+FN}$, where TP is the number of true positives, FP denotes the number of false positives, and FN tallies false negatives.

4.1 Experiments on Synthetic Data

We evaluated the four techniques over 5 synthetic binary image sets. These binary images were corrupted by zero mean Gaussian noise with unit standard deviation, and the task was to label the foreground objects (see the first and second columns in Fig. 1). Two of the sets contained balanced class labels (. . . and . . .), while the other three contained imbalanced classes. The five 150 image sets were divided into 100 images for training and 50 for testing. Example results and aggregate scores are shown in Fig. 1. Note that the last 4 columns illustrate the outcomes from each technique– SVMs, Logistic Regression (LR), SVRFs, and DRFs.

Logistic Regression and subsequently DRFs performed poorly in all three imbalanced data sets (. . . , . . . , and . . .). In these cases, SVMs outperformed these methods and consequently our proposed SVRFs outperformed

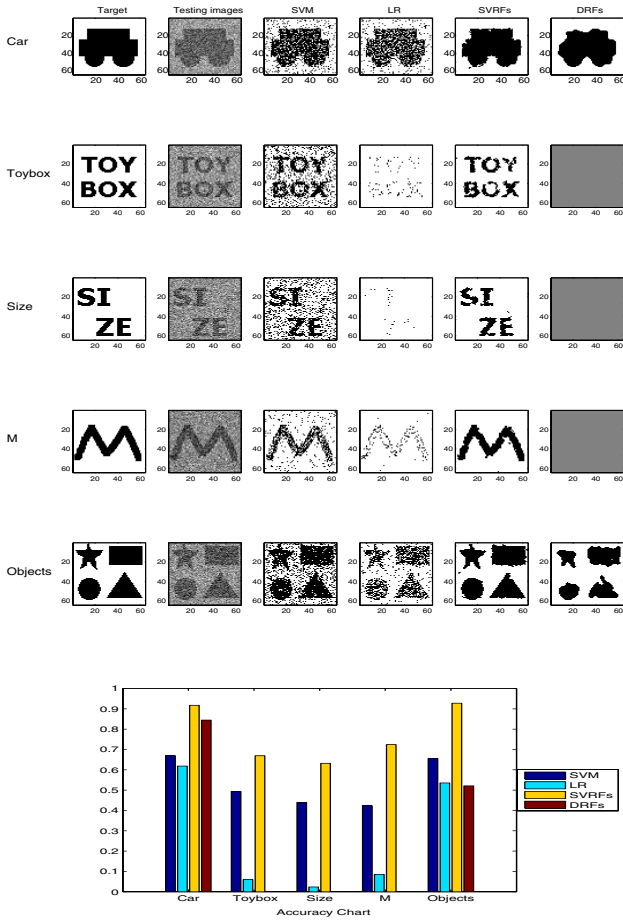


Fig. 1. Average scores on synthetic data sets

SVMs. In the first balanced data set (), DRFs and SVRFs both significantly outperformed SVMs and Logistic Regression (the iid classifiers). However, DRFs performed poorly on the second balanced data set (). This is due to DRFs simultaneous parameter learning, that tends to overestimate the local-consistency potential. Since the observation-matching is underweighted, edges become degraded during inference (there are more edge areas in the . . . data). Terminating inference before convergence could reduce this, but this is not highly desirable for automatic classification. Overall, our Support Vector Random Field model demonstrated the best performance on all data sets, in particular those with imbalanced data and a greater proportion of edge areas.

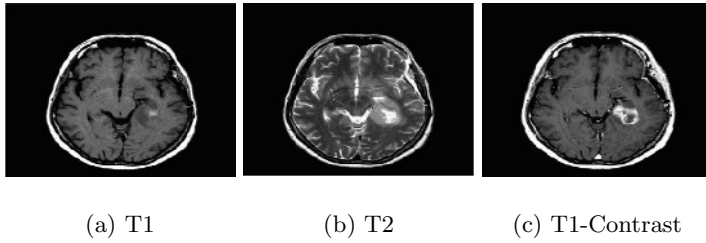


Fig. 2. A multi-spectral MRI

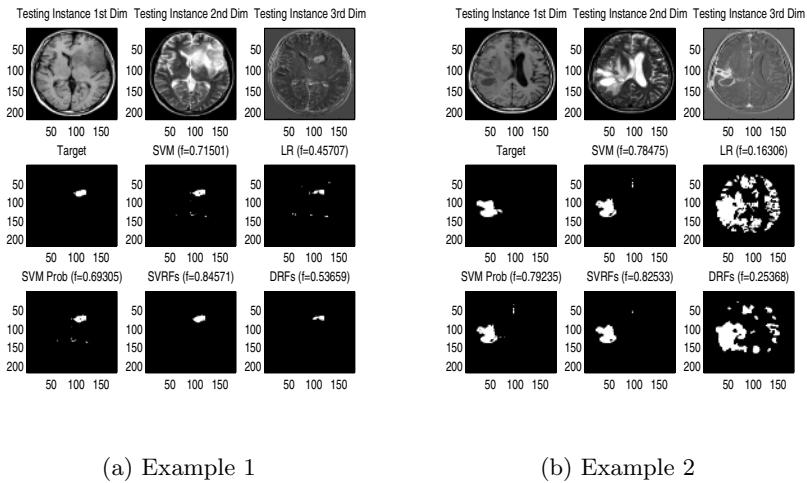


Fig. 3. An example of the classification result

4.2 Experiments on Real Data

We applied our model to the real-world problem of tumor segmentation in medical imaging. We focused on the task of brain tumor segmentation in MRI, an important task in surgical planning and radiation therapy currently being laboriously done by human medical experts. There has been significant research focusing on automating this challenging task (see [10]). Markov Random Fields have been explored previously for this task (see [10]), but recently SVMs have shown impressive performance [11,12]. This represents a scenario where our proposed Support Vector Random Field model could have a major impact. We evaluated the four classifiers from the previous section over 7 brain tumor patients. For each patient, three MRI ‘modalities’ were available: T1 (visualizing fat locations), T2 (visualizing water locations), and an additional T1 image with a ‘contrast agent’ added to enhance the visualization of metabolically active tumor areas (refer to Fig. 2).

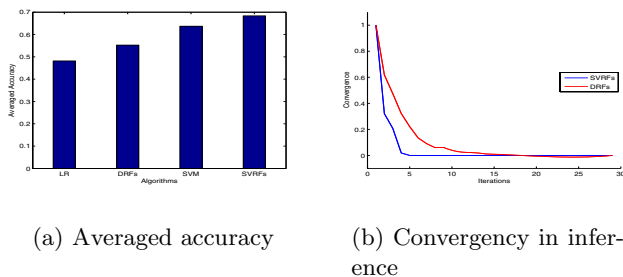


Fig. 4. Averaged accuracy and convergence in inference

The data was preprocessed with the Statistical Parametric Mapping software [13] to non-linearly align the images with a template in a standard coordinate system, and remove intensity inhomogeneity field effects. This non-linear template alignment approach was quantified to be highly effective in [14], and the inhomogeneity correction step computes a smooth corrective field that seeks to minimize the residual entropy after transformation of the log-intensity value’s probability distribution [15]. We used 12 features that incorporate image information and domain knowledge (the raw intensities, spatial expected intensities within the coordinate system, spatial priors for the brain area and normal tissue types within the coordinate system, the template image information, and left-to-right symmetry), each measured as features at 3 scales by using 3 different sizes of Gaussian kernel filters. We used a ‘patient-specific’ training scenario similar to [11,12].

Results for two of the patients are shown in Fig. 3, while average scores over the 7 patients are shown in Fig. 4(a). Note that ‘SVM+prob’ in Fig. 3 denotes the classification results from the Support Vector Machine posterior probability estimate. The Logistic Regression model performs poorly at this task, but DRFs perform significantly better. As with the synthetic data in cases of class imbalance, SVMs outperform both Logistic Regression and the DRFs. Finally, SVRFs improve the scores obtained by the SVMs by almost 5% (a significant improvement).

We compared convergence of the DRFs and SVRFs by measuring how many label changes occurred between inference iterations averaged over 21 trials (Fig. 4(a)). These results show that DRFs on average require almost 3 times as many iterations to converge, due to the overestimation of the local-consistency potential.

5 Conclusion

We have proposed a novel model for classification of data with spatial dependencies. The Support Vector Random Field combines ideas from SVMs and CRFs, and outperforms SVMs and DRFs on both synthetic data sets and an important real-world application. We also proposed an improvement to computing posterior

probability distributions from SVM decision functions, and a method to encourage continuity with local-consistency potentials. Our Support Vector Random Field model is robust to class imbalance, can be efficiently trained, converges quickly during inference, and can trivially be augmented with kernel functions to further improve results.

Acknowledgment

R. Greiner is supported by the National Science and Engineering Research Council of Canada (NSERC) and the Alberta Ingenuity Centre for Machine Learning (AICML). C.H. Lee is supported by NSERC, AICML, and iCORE. Our thanks to Dale Schuurmans for helpful discussions on optimization and parameter estimation, J. Sander for helpful discussions for the classification issues, BTGP members for help in data processing, and Albert Murtha (M.D.) for domain knowledge on the tumor data set.

References

1. Lafferty, J., Pereira, F., McCallum, A.: Conditional random fields: Probabilistic models for segmenting and labeling sequence data. *ICML* (2001)
2. Shawe-Taylor, Cristianini: *Kernel Methods for Pattern Analysis*. Cambridge University Press, Cambridge, UK (2004)
3. Li, S.Z.: *Markov Random Field Modeling in Image Analysis*. Springer-Verlag, Tokyo (2001)
4. Kumar, S., Hebert, M.: Discriminative random fields: A discriminative framework for contextual interaction in classification. *ICCV* (2003) 1150–1157
5. Besag, J.: On the statistical analysis of dirty pictures. *Journal of Royal Statistical Society. Series B* **48** (1986) 3:259–302
6. Kumar, S., Hebert, M.: Discriminative fields for modeling spatial dependencies in natural images. *NIPS* (2003)
7. R.Fletcher: *Practical Methods of Optimization*. John Wiley & Sons (1987)
8. Platt, J.: *Probabilistic outputs for support vector machines and comparison to regularized likelihood methods*. MIT Press, Cambridge, MA (2000)
9. Joachims, T.: Making large-scale svm learning practical. In Scholkopf, B., Burges, C., Smola, A., eds.: *Advances in Kernel Methods - Support Vector Learning*, MIT Press (1999)
10. Gering, D.: *Recognizing Deviations from Normalcy for Brain Tumor Segmentation*. PhD thesis, MIT (2003)
11. Zhang, J., Ma, K., Er, M., Chong, V.: Tumor segmentation from magnetic resonance imaging by learning via one-class support vector machine. *Int. Workshop on Advanced Image Technology* (2004) 207–211
12. Garcia, C., Moreno, J.: Kernel based method for segmentation and modeling of magnetic resonance images. *LNCS* **3315** (2004) 636–645
13. : Statistical parametric mapping, <http://www.fil.ion.bpmf.ac.uk/spm/> (Online)
14. Hellier, P., Ashburner, J., Corouge, I., Barillot, C., Friston, K.: Inter subject registration of functional and anatomical data using spm. In: *MICCAI*. Volume 587-590. (2002)
15. Ashburner, J.: Another mri bias correction approach. In: *8th Int. Conf. on Functional Mapping of the Human Brain*, Sendai, Japan. (2002)

Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication*

Jelenc Levec¹, Deepayan Chakrabarti¹,
Jon Kleinberg², and Christos Faloutsos¹

¹ School of Computer Science, Carnegie Mellon University
{jure, deepay, christos}@cs.cmu.edu

² Department of Computer Science, Cornell University
kleinber@cs.cornell.edu

Abstract. How can we generate realistic graphs? In addition, how can we do so with a mathematically tractable model that makes it feasible to analyze their properties rigorously? Real graphs obey a long list of surprising properties: Heavy tails for the in- and out-degree distribution; heavy tails for the eigenvalues and eigenvectors; small diameters; and the recently discovered “Densification Power Law” (DPL). All published graph generators either fail to match several of the above properties, are very complicated to analyze mathematically, or both. Here we propose a graph generator that is mathematically tractable and matches this collection of properties. The main idea is to use a non-standard matrix operation, the *Kronecker product*, to generate graphs that we refer to as “Kronecker graphs”.

We show that Kronecker graphs naturally obey all the above properties; in fact, we can rigorously *prove* that they do so. We also provide empirical evidence showing that they can mimic very well several real graphs.

1 Introduction

What do real graphs look like? How do they evolve over time? How can we generate synthetic, realistic, time-evolving graphs? Graph mining has been a fascinating challenge recently, with an emphasis on finding patterns and anomalies in social networks, communication networks, gene

* Work partially supported by the National Science Foundation under Grants No. IIS-0209107, SENSOR-0329549, IIS-0326322, CNS-0433540, CCF-0325453, IIS-0329064, CNS-0403340, CCR-0122581, a David and Lucile Packard Foundation Fellowship, and also by the Pennsylvania Infrastructure Technology Alliance (PITA), a partnership of Carnegie Mellon, Lehigh University and the Commonwealth of Pennsylvania’s Department of Community and Economic Development (DCED). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation, or other funding parties.

eg la o y ne wo , and any o e. Mo of he wo foc e on a ic na ho of g a h , whe e fa cina ing “law .” have been di cove ed, incl ding all dia e e , and heavy- ailed deg ee di ib ion .

A eali ic g a h gene a o i i o , an fo a lea wo ea on . The i ha i can gene a e g a h fo , ex a ola ion , “wha -if” cena io , and i lla ion , when eal g a h a e di c l o i o i ble o collec . Fo , exa ple , how well will a given oocol n on he In e ne ve yea f o now? Acc a e g a h gene a o can od ce o e eali ic odel fo he f e In e ne , on which i lla ion can be n. The econd ea on i o e b le: i fo ce o hin abo he a e n ha a g a h gene a o , ho ld obey , o be eali ic.

The ain con ib ion of hi a e a e he following:

- We ovide a gene a o , which obey all he ain a ic a e n ha have a ea ed in he li e a e .
- Gene a o al o obey he e cen ly di cove ed e o al evol ion a e n .
- Con a y o o he gene a o , ha a ch hi co bina ion of o e ie , o gene a o , lead o ac able analy i and o o oof .

O , gene a o , i ba ed on a non- anda d a ix o e a ion , he , The e a e eve al he o e on K onec e , od ce , which ac ally co e ond exac ly o a igni can o ion of wha we wan o ove: heavy- ailed di ib ion fo in-deg ee, o -deg ee, eigenval e , and eigenvec o . We al o de on a e how a K onec e G a h can a ch he behavio of eve al eal g a h (a en ci a ion , a e , ci a ion , and o he .). While K onec e , od ce have been d ied by he algeaic co bina o ic co ni y (ee e.g. [10]), he e en wo i he o e o e loy hi o e a ion in he de ign of ne wo odel o a ch eal da a e .

The e of he a e i o gani ed a follow : Sec ion 2 vey he e la ed li e a e . Sec ion 3 give he o o ed e hod. We e en he ex e i en al e l in Sec ion 4, and we clo e wi h o e di c ion and concl ion .

2 Related Work

Fi , we will di c the co only fo nd (a ic) a e n in g a h , hen o e e cen a e n on e o al evol ion, and nally, he a e of he a in g a h gene a ion e hod .

Static Graph Patterns: While any a e n have been di cove ed, wo of he inci al one a e heavy- ailed deg ee di ib ion and all dia e e .

The deg ee-di ib ion of a g a h i a owe law if he n be of node c_k wi h deg ee k i given by $c_k \propto k^{-\gamma}$ ($\gamma > 0$) whe e γ i called he owe-law ex onen . Powe law have been fo nd in he In e ne [13], he Web [15,7], ci a ion g a h [24], online ocial ne wo [9] and any o he . Devia ion fo he owe-law a e n have been no iced [23], which can be ex lained by he “DGX” di ib ion [5]. DGX i clo e ly e la ed o a nca ed logno al di ib ion.

Moreover, real-world graphs exhibit, relatively, all diameters (the “small-world” phenomenon): A graph has diameter d if every pair of nodes can be connected by a path of length at most d . The diameter d is noticeable only if $d \ll N$. Thus, a notable feature of the pairwise distance between nodes of a graph is the small-world property [26]. This is demonstrated by the initial number of hops in which communication (or transmission, say $q = 90\%$) of all connected pairs of nodes can reach each other. The effective diameter has been found to be small for large real-world graphs, like Inet, e-mail, Web, and social networks [2,21].

This is a result of the eigenvalue (or singular value) of the adjacency matrix of the graph, which is, in general, following a log-log scale. The characteristic value of the eigenvalue distribution obeys a power law. The distribution of eigenvalues (or singular values) (indicators of “network value”) has also been found to be skewed [9].

A number of other, even older, features have been found, including the “small-world” [14,9], “efficiency” [2,22], “clustering coefficient” and “any other”.

Temporal evolution Laws: Densification and shrinking diameter: Two very recent discoveries, both regarding the evolving graph, are worth mentioning [18]: (a) the “effective diameter” of a graph tends to decrease in overall size as the graph grows with time, and (b) the number of edges $E(t)$ and nodes $N(t)$ tend to obey the double power law (DPL), which can be written as

$$E(t) \propto N(t)^a \tag{1}$$

The parameter a is typically greater than 1, implying that the average degree of a node in the graph increases over time. This can be explained by a graph tending to have any node connected to more than one edge, and hence, defining a heavy growth.

Graph Generators: The easily obtainable generative model for a graph was a random graph model, where each pair of nodes has an identical, independent probability of being joined by an edge [11]. The study of this model has led to a rich analytical theory; however, this generative model for a graph has failed to capture real-world networks in a number of respects (for example, it does not capture heavy-tailed degree distribution).

The various types of recent models involve, of course, a number of features [1,2,28,15,16]: new nodes join the graph at each time step, and preferentially connect to existing nodes with high degree (the “rich get richer”). This simple behavior leads to power-law tails and to low diameters. The diameter in this model grows slowly with the number of nodes N , which violates the “shrinking diameter” property mentioned above.

Another family of graph-generating methods include, for all diameters, like the generative model [27] and the Waxman generative model [6]. A hybrid family of methods how heavy-tailed degree if nodes are only interested in connecting with nodes of degree constraint [8,12].

Summary: Most current generative models focus on only one (or a few) features, and neglect the other. In addition, it is generally hard to overcome some of the features. The generative models described in the next section address these issues.

3 Proposed Method

The method we propose is based on a recursive construction. Denoting the recursive function solely in words what is possible, a number of standard graph construction methods fail to detect graph hardness according to the algorithm observed in practice, and they also do not catch who is doing the increase. To detect denoting graph with construction, and hereby, a change in the overall behavior of real networks, we develop a procedure that is better described in terms of the recursive nature of practice. To help in the definition of the method, the accompanying table provides a list of symbols and their definition.

Sy mbol	De finition
G_1	the initialization of a Kronecker Graph
N_1	number of nodes in initialization
E_1	number of edges in initialization
$G_1^{[k]} = G_k$	the k^{th} Kronecker power of G_1
a	denotation expansion
d	degree of a graph
\mathcal{P}_1	probability matrix

3.1 Main Idea

The main idea is to create self-similar graphs, recursively. We begin with an initial graph G_1 , with N_1 nodes and E_1 edges, and by recursion we produce successively larger graphs $G_2 \dots G_n$. Each time the k^{th} graph G_k is on $N_k = N_1^k$ nodes. If we want the graph to exhibit a violation of the Denotation Power Law, then G_k should have $E_k = E_1^k$ edges. This is also why having a degree cascade in order to get a standard recursive construction (for example, the additional Cartesian product of the construction of [4]) do not satisfy it.

In order to have the recursive nature of work, a change in the effectiveness of this goal. The Kronecker product is defined as follows:

Definition 1 (Kronecker product of matrices). Let $A = [a_{i,j}]$ be an $n \times m$ matrix and $B = [b_{i',j'}]$ be an $n' \times m'$ matrix. Then the Kronecker product of A and B is the $(n * n') \times (m * m')$ matrix $C = A \otimes B$ defined as:

$$C = A \otimes B \doteq \begin{pmatrix} a_{1,1}B & a_{1,2}B & \dots & a_{1,m}B \\ a_{2,1}B & a_{2,2}B & \dots & a_{2,m}B \\ \dots & \dots & \dots & \dots \\ a_{n,1}B & a_{n,2}B & \dots & a_{n,m}B \end{pmatrix} \tag{2}$$

We define the Kronecker product of two graphs as the Kronecker product of their adjacency matrices.

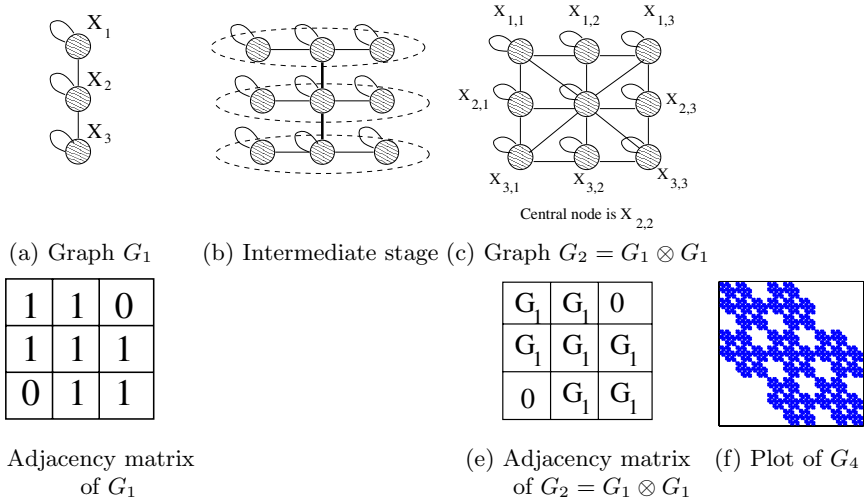


Fig. 1. Example of Kronecker multiplication: Top: a “3-chain” and its Kronecker product with itself; each of the X_i nodes gets expanded into 3 nodes, which are then linked using Observation 1. Bottom row: the corresponding adjacency matrices, along with matrix for the fourth Kronecker power G_4 .

Observation 1 (Edges in Kronecker-multiplied graphs)

$$(X_{ij}, X_{kl}) \in G \otimes H \iff (X_i, X_k) \in G \text{ and } (X_j, X_l) \in H$$

The last observation is a subtle, but crucial, and deserve elaboration: Figure 1(a–c) shows the recursive construction of $G \otimes H$, when $G = H$ is a 3-node graph. Consider node $X_{1,2}$ in Figure 1(c): It belongs to the H graph having placed node X_1 (see Figure 1(b)), and in fact it is the X_2 node (i.e., the center) within this all H -graph.

We now observe a growing influence of graph by iterating the Kronecker product:

Definition 2 (Kronecker power). The k^{th} Kronecker power of graph G_1 is denoted by $G_1^{[k]}$ (where $G_1^{[1]} = G_1$).

$$G_1^{[k]} = G_k = \underbrace{G_1 \otimes G_1 \otimes \dots \otimes G_1}_{k \text{ times}} = G_{k-1} \otimes G_1$$

The self-similarity of the Kronecker graph is clear: To produce G_k from G_{k-1} , we “expand” (i.e. place) each node of G_{k-1} by connecting it into a copy of G , and we join the edge together, according to the adjacency in

G_{k-1} (see Fig. 1). This process is very natural: one can imagine a growing hierarchical community with the growth, respectively, with nodes in the community, respectively getting expanded in original community of the community. Node in the community then links along the level and all nodes form different community.

3.2 Theorems and Proofs

We shall now discuss the properties of Konec degree hierarchy, specifically, hierarchical distribution, diameter, eigenvalue, eigenvector, and time-evolution. Our ability to solve analytical problems about all of the properties is a major advantage of Konec degree hierarchy over other generalizations. The next few theorems solve hierarchical distribution of inner area, linear form of Konec degree hierarchy model. This is important, because a careful choice of the initial graph G_1 can make the hierarchical distribution to behave like a power-law or DGX distribution.

Theorem 1 (Multinomial degree distribution). *Let G_1 be a graph with N_1 nodes and E_1 edges. Let d_1, d_2, \dots, d_{N_1} be the degrees of the nodes in G_1 . Let G_k be the k -th iteration of the Konec degree hierarchy construction. Then the degree distribution of G_k is given by*

Let the initial G_1 have the degree sequence d_1, d_2, \dots, d_{N_1} . Konec degree hierarchy of a node with degree d expand into N_1 nodes, with the corresponding degree being $d \times d_1, d \times d_2, \dots, d \times d_{N_1}$. After Konec degree hierarchy, the degree of each node in graph G_k is of the form $d_{i_1} \times d_{i_2} \times \dots \times d_{i_k}$, with $i_1, i_2, \dots, i_k \in \{1, \dots, N_1\}$, and here is one node for each ordered combination. This gives the multinomial distribution on the degree of G_k . No other has the degree of node in G_k can be expressed as the k^{th} Konec degree of the vector $(d_1, d_2, \dots, d_{N_1})$. □

Theorem 2 (Multinomial eigenvalue distribution). *Let G_1 be a graph with N_1 nodes and E_1 edges. Let $\lambda_1, \lambda_2, \dots, \lambda_{N_1}$ be the eigenvalues of G_1 . Let G_k be the k -th iteration of the Konec degree hierarchy construction. Then the eigenvalue distribution of G_k is given by*

Let G_1 have the eigenvalue $\lambda_1, \lambda_2, \dots, \lambda_{N_1}$. By properties of the Konec degree hierarchy [19,17], the eigenvalue of G_k are k^{th} Konec degree of the vector $(\lambda_1, \lambda_2, \dots, \lambda_{N_1})$. As in Theorem 1, the eigenvalue distribution is a multinomial. □

An illustrative example of the properties of Konec degree hierarchy is given below. How the following.

Theorem 3 (Multinomial eigenvector distribution). *Let G_1 be a graph with N_1 nodes and E_1 edges. Let v_1, v_2, \dots, v_{N_1} be the eigenvectors of G_1 . Let G_k be the k -th iteration of the Konec degree hierarchy construction. Then the eigenvector distribution of G_k is given by*

We have just covered every aspect of the hierarchical structure. Notice that the proof we did is consistent of the Konec degree hierarchy properties. Next we continue with the overall structure: the denotation, power law, and highlighting/ability of diameter.

Theorem 4 (DPL). $(G_k)_{k \geq 1}$ is a DPL with $a = \log(E_1)/\log(N_1)$.

Since the k^{th} Konec edge G_k has $N_k = N_1^k$ nodes and $E_k = E_1^k$ edges, it satisfies $E_k = N_k^a$, where $a = \log(E_1)/\log(N_1)$. The crucial point is that this exponent a is independent of k , and hence the sequence of Konec edges follows an exact version of the Denicaion Power Law. \square

We now show how the Konec edge model allows us to view the evolution of a network as a special ingredient for a changing diameter. The evolution of any real-world network is a process. In order to be able to handle it, we will assume that the initial graph G_1 has a self-loop on every node; otherwise, the Konec edges may in fact be disconnected.

Lemma 1. $G \otimes H$ has diameter d if and only if G and H have diameter d .

Each node in $G \otimes H$ can be seen as an ordered pair (v, w) , with v a node of G and w a node of H , and with an edge joining (v, w) and (x, y) exactly when (v, x) is an edge of G and (w, y) is an edge of H . Now, for an arbitrary pair of nodes (v, w) and (v', w') , we show how to find a path of length at most d connecting them. Since G has diameter at most d , there is a path $v = v_1, v_2, \dots, v_r = v'$, where $r \leq d$. If $r < d$, we can convert this into a path $v = v_1, v_2, \dots, v_d = v'$ of length exactly d , by simply repeating v' at the end for $d - r$ steps. By an analogous argument, we have a path $w = w_1, w_2, \dots, w_d = w'$. Now by the definition of the Konec edge model, there is an edge joining (v_i, w_i) and (v_{i+1}, w_{i+1}) for all $1 \leq i \leq d - 1$, and so $(v, w) = (v_1, w_1), (v_2, w_2), \dots, (v_d, w_d) = (v', w')$ is a path of length d connecting (v, w) to (v', w') , as desired. \square

Theorem 5. G_k has diameter d if and only if G has diameter d .

This follows directly from the previous lemma, combined with induction on k . \square

We also consider the effective diameter d_e ; we define the q -effective diameter as the minimum d_e such that, for a least a q fraction of the reachable nodes, the path length is at most d_e . The q -effective diameter is a non-obvious quantity to handle, because, while the latter being a measure of the effectiveness of degeneration, it is not clear in the graph (e.g. very long chains); however, the q -effective diameter and diameter tend to exhibit similar behavior. For example, in the binomial distribution, we will generally consider the q -effective diameter with $q = .9$, and, effectively, it is the same as the diameter.

Theorem 6 (Effective Diameter). G_k has effective diameter d if and only if G has effective diameter d .

To prove this, it is sufficient to show that for every node v of G_k , the probability that the distance d converges to a k -gon in n is $1 - (1 - 2/N_1)^k$.

We establish this as follows. Each node in G_k can be seen as an ordered sequence of k nodes from G_1 , and we can view the random selection of a node in G_k as a sequence of k independent random node selections from G_1 . Suppose we have $v = (v_1, \dots, v_k)$ and $w = (w_1, \dots, w_k)$ as two randomly selected nodes from G_k . Now, if x and y are two nodes in G_1 at distance d (such a pair (x, y) exists since G_1 has diameter d), then with probability $1 - (1 - 2/N_1)^k$, the i th node in v is x and the j th node in w is y . If the i th and j th nodes are x and y , then the distance between v and w is d . As the expected value of $1 - (1 - 2/N_1)^k$ converges to a k -gon, it follows that the q -effective diameter is converging to d . \square

3.3 Stochastic Kronecker Graphs

While the Kronecker power construction did not yield graphs with a range of desired properties, it did create a natural “ai-ca-e effect” in the degree and local structure, primarily because individual values have large influence. Here we propose a stochastic version of Kronecker graphs that eliminates this effect completely.

We start with an $N_1 \times N_1$ matrix \mathcal{P}_1 : the value p_{ij} denotes the probability that edge (i, j) is selected. We construct the k th Kronecker power $\mathcal{P}_1^{[k]} = \mathcal{P}_k$; and then for each entry p_{uv} of \mathcal{P}_k , we include an edge between node u and v with probability $p_{u,v}$. The resulting binary adjacency matrix $R = R(\mathcal{P}_k)$ will be called the k th Kronecker power of \mathcal{P}_1 (or k -Kronecker power).

In principle one could simply choose each of the N_1^2 entries of the matrix \mathcal{P}_1 independently. However, we need to be careful about a few things: α and β . Let G_1 be the initial adjacency matrix (binary, degree α and β); we create the corresponding probability matrix \mathcal{P}_1 by replacing each “1” and “0” of G_1 with α and β , respectively ($\beta \leq \alpha$). The resulting probability matrix is again — with some adjustments — the self-similar structure of the Kronecker graph in the previous section (which, for clarity, we call k -Kronecker power).

We need to carefully handle the graph model used by this model construction to exhibit the desired properties of real data sets, and without the ai-ca-e effect of the degree distribution. The goal of setting α and β on a chosen observed data set is to control the degree distribution, to provide the core of the data set. In our experiments in the construction, we use the metric which we describe here.

4 Experiments

Now, we demonstrate the ability of Kronecker graphs to match the structure of real-world graphs. The data sets we use are:

- **Facebook**: This is a citation graph for high-energy physics papers, with a total of $N = 29,555$ nodes and $E = 352,807$ citations. We follow the evolution from January 1993 to April 2003, with one data point every month.

• This is a U.S. encyclopedia data set that has an average of Jan a, y 1963 to Dece mber, 1999. The graph contains a total of $N = 3,942,825$ nodes and $E = 16,518,948$ edges. Citations graph are naturally considered a directed graph. For the purpose of this work we have undirected.

• We also analyze a data set consisting of a single national highway network. In the network only edges for Jan a, y 2000, with $N = 6,474$ and $E = 26,467$.

We observe two kinds of graphs are seen — “basic” and “evolutional.” A common example, common to basic data sets include the degree distribution, the characteristic (eigenvalue of graph adjacency matrix, λ), principal eigenvector of adjacency matrix and the distribution of connected components. The evolutional data sets include the diameter, average, height of the giant component, average, and the density of the network. For the diameter component, we have a good understanding of the effective diameter, has a relatively small number of nodes and effective diameter, but the linear evolution of a set of nodes on non-eigenvalue; see [18] for further details on this calculation.

Results are shown in Figure 2 and 3 for the graphs which evolve over time (basic and evolutional). For brevity, we show the results for only two basic and two evolutional data sets. We see that the degree, initial Konec, model, early characteristics, the average degree and eigenvalue distribution, as well as the evolutional data sets are defined by the Density Power Law and the ability of diameter. However, the degree, initial network model, early characteristics, the average degree, as shown in characteristic for the degree, initial Konec, graph of Figure 2 (second row, second column). We see that the Stochastic Konec, graph, the evolution of the distribution, further, showing the average degree, characteristic of the real data; they all show a characteristic in the before-ability of the end of the diameter of real graphs.

For the Stochastic Konec, graph we need to define the parameters α and β defined in the previous section. This leads to interesting question who are the outliers lie beyond the core of the network; consequently, we reached by the force over (the relatively small number of) possible initial graphs of the network, and we then choose α and β to approach well the edge density, the average degree, the spectral properties, and the DPL exponent.

Finally, Figure 4 shows the results for the basic data sets in the evolutional graphs. Recall that we analyze a single, basic network in this case. In addition to the degree distribution and characteristic, we also show two typical results [9]: the distribution of principal eigenvector components, λ , (eigenvalue, λ) and the network (the number of eachable sites $P(h)$ within h hops of a node), as a function of the number of hops (h).

5 Observations and Conclusions

Here we will review of the desirable properties of the evolutional Konec, graph and Stochastic Konec, graph.

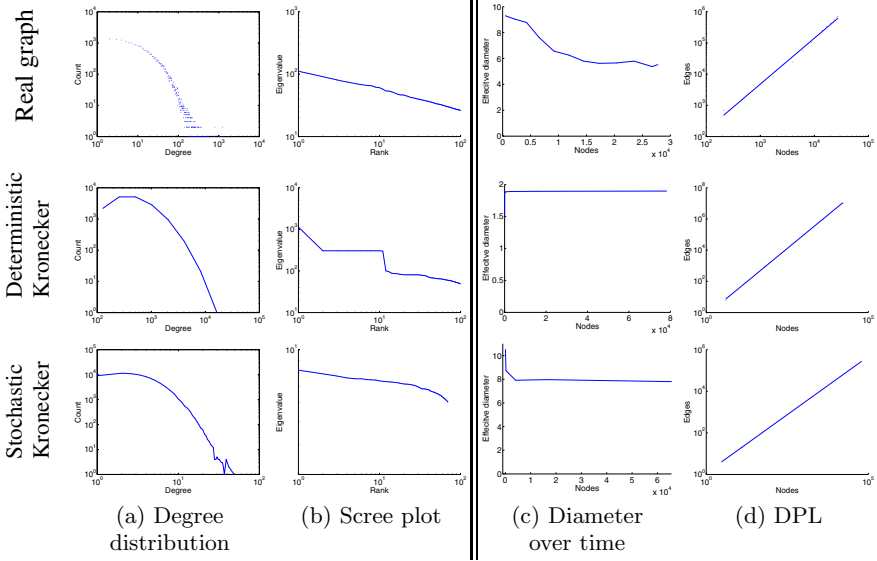


Fig. 2. *arXiv* dataset: Patterns from the real graph (top row), the deterministic Kronecker graph with G_1 being a star graph with 3 satellites (middle row), and the Stochastic Kronecker graph ($\alpha = 0.41$, $\beta = 0.11$ – bottom row). *Static* patterns: (a) is the PDF of degrees in the graph (log-log scale), and (b) the distribution of eigenvalues (log-log scale). *Temporal* patterns: (c) gives the effective diameter over time (linear-linear scale), and (d) is the number of edges versus number of nodes over time (log-log scale). Notice that the Stochastic Kronecker Graph qualitatively matches all the patterns very well.

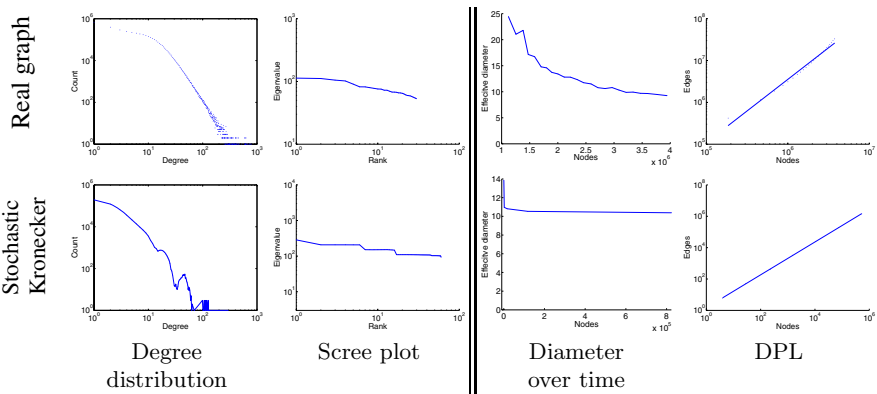


Fig. 3. *Patents*: Again, Kronecker graphs match all of these patterns. We show only the Stochastic Kronecker graph for brevity.

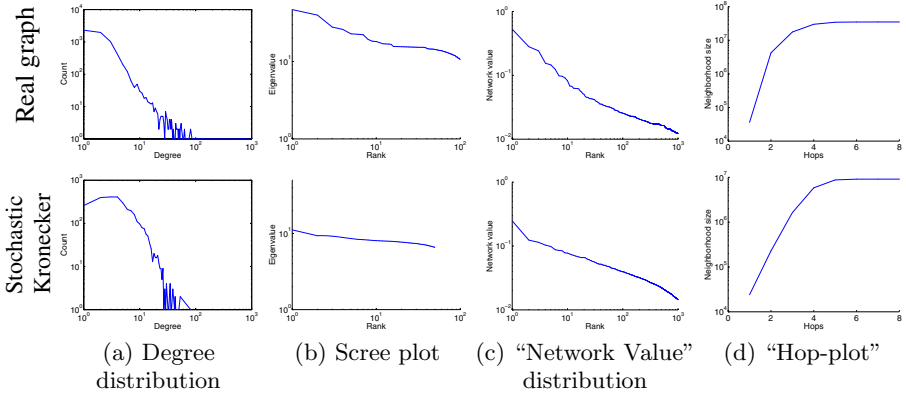


Fig. 4. *Autonomous systems:* Real (top) versus Kronecker (bottom). Columns (a) and (b) show the degree distribution and the scree plot, as before. Columns (c) and (d) show two more static patterns (see text). Notice that, again, the Stochastic Kronecker Graph matches well the properties of the real graph.

Stochastic Kronecker Graph include eigenvalue generation, a special case: For $\alpha=\beta$, we obtain an Erdős-Renyi graph; for $\alpha=1$ and $\beta=0$, we obtain a deterministic Kronecker graph; using the G_1 matrix of a 2×2 matrix, we obtain the RMAT generation [9]. In contrast to Kronecker graph, the RMAT cannot exist alone in itself, since it needs to know the number of edges in it. Thus, it is incapable of obeying the “density law”.

The Erdős-Renyi graph exhibits phase transition [11]. Several real world graphs have a “edge of chaos” [3,25]. In contrast to Stochastic Kronecker Graph also exhibits phase transition. For all values of α and β , Stochastic Kronecker Graph have any all disconnected components; for large values they have a giant component with all diameters. In between, they exhibit behavior suggestive of a phase transition: For a carefully chosen (α, β) , the diameters, large, and a giant component just appear, emerging. We will see details, for lack of space.

All of these are false for the deterministic Kronecker Graph. However, the idea of working on the properties of and analysis (see e.g. [20]), which one could potentially apply in order to overcome the limitations of the Stochastic Kronecker Graph.

In conclusion, the main contribution of this work is a family of graph generation, using a non-traditional matrix operation, the Kronecker product. The resulting graphs (a) have all the characteristic (heavy-tailed degree distribution, all diameters), (b) all the essential properties (density law, thinning diameters), and in addition, (c) we can formally prove all of the properties.

Several of the proofs are extremely simple, thanks to the richness of Kronecker multiplication. We also provide proofs about the diameters and “effective diameters”, and we show how the Stochastic Kronecker Graph can be used to simulate real graphs well.

References

1. R. Albert and A.-L. Barabasi. Emergence of scaling in random networks. *Science*, pages 509512, 1999.
2. R. Albert and A.-L. Barabasi. Statistical mechanics of complex networks. *Reviews of Modern Physics*, 2002.
3. P. Bak. How nature works : The science of self-organized criticality, Sept. 1996.
4. A.-L. Barabasi, E. Ravasz, and T. Vicsek. Deterministic scale-free networks. *Physica A*, 299:559-564, 2001.
5. Z. Bi, C. Faloutsos, and F. Korn. The DGX distribution for mining massive, skewed data. In *KDD*, pages 17-26, 2001.
6. B.M.Waxman. Routing of multipoint connections. *IEEE Journal on Selected Areas in Communications*, 6(9), December 1988.
7. A. Broder, R. Kumar, F. Maghoul1, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Wiener. Graph structure in the web: experiments and models. In *Proceedings of World Wide Web Conference*, 2000.
8. J. M. Carlson and J. Doyle. Highly optimized tolerance: a mechanism for power laws in designed systems. *Physics Review E*, 60(2):1412-1427, 1999.
9. D. Chakrabarti, Y. Zhan, and C. Faloutsos. R-MAT: A recursive model for graph mining. In *SIAM Data Mining*, 2004.
10. T. Chow. The Q-spectrum and spanning trees of tensor products of bipartite graphs. *Proc. Amer. Math. Soc.*, 125:3155-3161, 1997.
11. P. Erdos and A. Renyi. On the evolution of random graphs. Publication of the Mathematical Institute of the Hungarian Academy of Science, 5:17-67, 1960.
12. A. Fabrikant, E. Koutsoupias, and C. H. Papadimitriou. Heuristically optimized trade-offs: A new paradigm for power laws in the internet (extended abstract), 2002.
13. M. Faloutsos, P. Faloutsos, and C. Faloutsos. On power-law relationships of the internet topology. In *SIGCOMM*, pages 251-262, 1999.
14. M. Girvan and M. E. J. Newman. Community structure in social and biological networks. In *Proc. Natl. Acad. Sci. USA*, volume 99, 2002.
15. J. M. Kleinberg, S. R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. The web as a graph: Measurements, models and methods. In *Proceedings of the International Conference on Combinatorics and Computing*, 1999.
16. S. R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. Extracting large-scale knowledge bases from the web. *VLDB*, pages 639-650, 1999.
17. A. N. Langville and W. J. Stewart. The Kronecker product and stochastic automata networks. *Journal of Computation and Applied Mathematics*, 167:429-447, 2004.
18. J. Leskovec, J. Kleinberg, and C. Faloutsos. Graphs over time: densification laws, shrinking diameters and possible explanations. In *KDD'05*, Chicago, IL, USA, 2005.
19. C. F. V. Loan. The ubiquitous Kronecker product. *Journal of Computation and Applied Mathematics*, 123:85-100, 2000.
20. M. Mehta. Random Matrices. Academic Press, 2nd edition, 1991.
21. S. Milgram. The small-world problem. *Psychology Today*, 2:60-67, 1967.
22. C. R. Palmer, P. B. Gibbons, and C. Faloutsos. Anf: A fast and scalable tool for data mining in massive graphs. In *SIGKDD*, Edmonton, AB, Canada, 2002.
23. D. M. Pennock, G.W. Flake, S. Lawrence, E. J. Glover, and C. L. Giles. Winners dont take all: Characterizing the competition for links on the web. *Proceedings of the National Academy of Sciences*, 99(8):5207-5211, 2002.

24. S. Redner. How popular is your paper? an empirical study of the citation distribution. *European Physical Journal B*, 4:131-134, 1998.
25. R. Sole and B. Goodwin. *Signs of Life: How Complexity Pervades Biology*. Perseus Books Group, New York, NY, 2000.
26. S. L. Tauro, C. Palmer, G. Siganos, and M. Faloutsos. A simple conceptual model for the internet topology. In *Global Internet, San Antonio, Texas*, 2001.
27. D. J. Watts and S. H. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393:440-442, 1998.
28. J. Winick and S. Jamin. Inet-3.0: Internet Topology Generator. Technical Report CSE-TR-456-02, University of Michigan, Ann Arbor, 2002.

A Correspondence Between Maximal Complete Bipartite Subgraphs and Closed Patterns

Jinyan Li, Haiyan Li, Donny Soh, and Limsoon Wong

Institute for Infocomm Research,
21 Heng Mui Keng Terrace, Singapore 119613
{jinyan, haiquan, studonny, limsoon}@i2r.a-star.edu.sg

Abstract. For an undirected graph G without self-loop, we prove: (i) that the number of closed patterns in the adjacency matrix of G is even; (ii) that the number of the closed patterns is precisely double the number of maximal complete bipartite subgraphs of G ; (iii) that for every maximal complete bipartite subgraph, there always exists a unique pair of closed patterns that matches the two vertex sets of the subgraph. Therefore, we can enumerate all maximal complete bipartite subgraphs by using efficient algorithms for mining closed patterns which have been extensively studied in the data mining field.

1 Introduction

In recent years, graph and heuristic algorithms have grown to a very broad extent in the past decade (see [18] and the Preface of [8]), largely due to the effectiveness of graph models in many areas such as theoretical, electrical engineering, computer logic, business administration, sociology, economics, marketing, biology, and networking and communication. In addition, any problem can be modeled with graphs, \dots , (see the definition below) formed by grouping two non-overlapping sets of vertices of a certain graph having all kinds of full connections between them.

We consider two examples. Suppose the telephone network in a mobile communication network. Some people have a wide range of contacts, while others have few. Which groups of contacts (with a maximal number) have a full interaction with another group of contacts? This question can be modeled by a graph where a mobile phone contact is a node and a communication is an edge. Then, a maximal bipartite subgraph of this graph corresponds to two groups of contacts between which there exists a full communication. This problem is similar to the one studied in web mining [4,9,12] where web communities are modeled by bipartite graphs. Another example is about 'protein' interaction in a cell. The naturally horizontal and vertical protein in a cell have interaction with one another. This question again can be modeled by a graph, where a protein is a node and an interaction between a pair of proteins forms an edge. Then, listing all maximal bipartite subgraphs of this graph can answer the question which two protein groups have a full interaction, which is a problem studied in biology [14,15].

Using all axial complete bipartite graphs has been studied theoretically in [5]. The problem of all axial complete bipartite graphs of a graph can be enumerated in time $O(a^3 2^{2a} n)$, where a is the adjacency of the graph and n is the number of vertices in the graph. Even though the algorithm has a linear complexity, it is not practical for large graphs due to the large constant overhead (a can easily be around 10-20 in practice) [20]. In this paper, we study this problem for a data mining exercise: We use heuristic data mining algorithms to efficiently enumerate all axial complete bipartite graphs for a large graph. A main concept of the data mining algorithm is called *axial complete bipartite graph*. The enumeration algorithm and its implementation devoted to the mining of closed itemsets for the so-called *axial complete bipartite graph* [2,6,7,13,16,17,19]. The data mining algorithm and the mining procedure are end-to-end. Our main contribution here is the observation that the mining of closed itemsets for the adjacency matrix of a graph, used as a special analytical database, is equivalent to the problem of enumerating all axial complete bipartite graphs of the graph.

The rest of this paper is organized as follows: Section 2 and 3 provide basic definitions and notation on graphs and closed itemsets. In Section 4 we prove that there is a one-to-one correspondence between closed itemsets and axial complete bipartite graphs for any finite graph. In Section 5, we present an experimental evaluation on a 'real world' transaction graph. Section 6 discusses related work and then concludes this paper.

2 Maximal Complete Bipartite Subgraphs

A **graph** $G = \langle V^G, E^G \rangle$ is composed of a set of vertices V^G and a set of edges $E^G \subseteq V^G \times V^G$. We often omit the superscripts in V^G, E^G and often place when the context is clear. Throughout this paper, we assume G is an undirected graph without any self-loops. In other words, we assume that (i) there is no edge $(u, u) \in E^G$ and (ii) for every $(u, v) \in E^G$, (u, v) can be replaced by (v, u) —that is, (u, v) is an undirected edge.

A graph H is a **subgraph** of a graph G if $V^H \subseteq V^G$ and $E^H \subseteq E^G$. A graph G is **bipartite** if V^G can be partitioned into two non-empty and non-intersecting subsets V_1 and V_2 such that $E^G \subseteq V_1 \times V_2$. This bipartite graph is usually denoted by $G = \langle V_1 \cup V_2, E^G \rangle$. Note that there is no edge in G joining two vertices within V_1 or V_2 . G is **complete bipartite** if $V_1 \times V_2 = E^G$.

Two vertices u, v of a graph G are said to be adjacent if $(u, v) \in E^G$ —that is, there is an edge in G that connects them. The **neighborhood** $\beta^G(v)$ of a vertex v of a graph G is the set of all vertices in G that are adjacent to v —that is, $\beta^G(v) = \{u \mid (u, v) \text{ or } (v, u) \in E^G\}$. The neighborhood $\beta^G(X)$ for a subset X of vertices of a graph G is the set of common neighborhoods of the vertices in X —that is, $\beta^G(X) = \bigcap_{x \in X} \beta^G(x)$.

Note that for any subset X of vertices of a graph G , such as X and $\beta^G(X)$ are both non-empty, it is clear that $H = \langle X \cup \beta^G(X), X \times \beta^G(X) \rangle$ is a complete bipartite graph of G . Note also that it is possible for a vertex $v \notin X$ of G

be adjacent to every vertex of $\beta^G(X)$. In this case, the set X can be extended by adding the vertex v , while maintaining the same neighborhood. Where do the extensions go? We use the following definition of maximal complete bipartite subgraph.

Definition 1. Let $H = \langle V_1 \cup V_2, E \rangle$ be a complete bipartite subgraph of G . H is called a **maximal complete bipartite subgraph** of G if $\beta^G(V_1) = V_2$ and $\beta^G(V_2) = V_1$.

Not all maximal complete bipartite subgraphs are equally interesting. Recall that each linear ordering of vertices induces a total ordering on the complete bipartite subgraphs. We would probably not be very interested in those with a small size. We would probably be considerably more interested if one of the two parts is large, or both of the two parts are large. Hence, we introduce the notion of density on maximal complete bipartite subgraphs.

Definition 2. Let $H = \langle V_1 \cup V_2, E \rangle$ be a complete bipartite subgraph of G . Let (m, n) be a pair of positive integers such that $|V_1| = m$ and $|V_2| = n$. The **density** of H with respect to (m, n) is defined as $\frac{|E|}{m \cdot n}$.

A complete bipartite subgraph $H = \langle V_1 \cup V_2, E \rangle$ of G is called a **maximal complete bipartite subgraph** of G if $\beta^G(V_1) = V_2$ and $\beta^G(V_2) = V_1$ and it is **maximal** if there is no other complete bipartite subgraph $H' = \langle V'_1 \cup V'_2, E' \rangle$ of G with $V_1 \subset V'_1$ and $V_2 \subset V'_2$ such that $\beta^G(V'_1) = V'_2$ and $\beta^G(V'_2) = V'_1$. To appreciate this notion of maximality, we prove the proposition below.

Proposition 1. Let $H = \langle V_1 \cup V_2, E \rangle$ and $H' = \langle V'_1 \cup V'_2, E' \rangle$ be complete bipartite subgraphs of G such that $V_1 \subseteq V'_1$ and $V_2 \subseteq V'_2$. If $H = H'$ and $\beta^G(V_1) = V_2$ and $\beta^G(V_2) = V_1$, then $\beta^G(V'_1) = V'_2$ and $\beta^G(V'_2) = V'_1$. In other words, if H is a maximal complete bipartite subgraph of G , then H' is also a maximal complete bipartite subgraph of G .

3 Closed Patterns of an Adjacency Matrix

The adjacency matrix of a graph is an interesting object. Let G be a graph with $V^G = \{v_1, v_2, \dots, v_p\}$. The **adjacency matrix** \mathbf{A} of G is the $p \times p$ matrix defined by

$$\mathbf{A}[i, j] = \begin{cases} 1 & \text{if } (v_i, v_j) \in E^G \\ 0 & \text{otherwise} \end{cases}$$

Recall that graphs do not have self-loops and are undirected. Thus \mathbf{A} is a symmetric matrix and every entry on the main diagonal is 0. Also, $\{v_j \mid \mathbf{A}[k, j] = 1, 1 \leq j \leq p\} = \beta^G(v_k) = \{v_j \mid \mathbf{A}[j, k] = 1, 1 \leq j \leq p\}$.

The adjacency matrix of a graph can be interpreted in a **transactional database** (DB) [1]. To define a DB , we first define a **transaction**. Let I be a set of **items**. Then a transaction is defined as a subset of I . For example, a set I is to be all items in a store, a transaction by a customer is the items he has bought. A DB is a non-empty list of transactions. Each transaction T in a DB is assigned a unique identity $id(T)$. A **pattern** is defined as a non-empty subset of items of I . A pattern may or may not be contained in a transaction. Given a DB and a pattern P , then the set of transactions in DB containing P is called the **support** of P , denoted $sup^{DB}(P)$. We are often interested in a pattern having sufficiently frequent in a DB . The pattern is called **frequent** pattern — having a pattern P satisfying $sup^{DB}(P) \geq ms$, for a threshold $ms > 0$. In this case, namely mentioned otherwise, we consider all and only those patterns with a non-zero support, namely all those frequent patterns with the support hold $ms = 1$. So, by a pattern of a DB , we mean having a non-empty and occurring in DB at least once.

Let G be a graph with $V^G = \{v_1, v_2, \dots, v_p\}$. If each vertex in V^G is defined as an item, then the neighborhood $\beta^G(v_i)$ of v_i is a transaction. Thus,

$$\{\beta^G(v_1), \beta^G(v_2), \dots, \beta^G(v_p)\}$$

is a DB . Such a special DB is denoted by DB_G . The identity of a transaction in DB_G is defined as the vertex itself — having, $id(\beta^G(v_i)) = v_i$. Note that DB_G has the pattern set of items and transactions. Note also that $v_i \notin \beta^G(v_i)$ since we assume G to be an undirected graph without self-loops.

DB_G can be represented as a binary adjacency matrix. This binary matrix \mathbf{B} is defined by

$$\mathbf{B}[i, j] = \begin{cases} 1 & \text{if } v_j \in \beta^G(v_i) \\ 0 & \text{otherwise} \end{cases}$$

Since $v_j \in \beta^G(v_i)$ iff $(v_i, v_j) \in E^G$, it can be seen that $\mathbf{A} = \mathbf{B}$. So, “a pattern of DB_G ” is equivalent to “a pattern of the adjacency matrix of G ”.

Closed patterns are a type of interesting pattern in a DB . In the last few years, the problem of efficiently mining closed patterns for a large DB has attracted a lot of researchers in the data mining community [2,6,7,13,16,17,19]. Let I be a set of items, and D be a transactional database defined on I . For a pattern $P \subseteq I$, let $f^D(P) = \{T \in D \mid P \subseteq T\}$ — having, $f^D(P)$ are all transactions in D containing the pattern P . For a set of transactions $D' \subseteq D$, let $g(D') = \bigcap_{T \in D'} T = \bigcap D'$ — having, the set of items which are shared by all transactions in D' . Using the above function, we can define the notion of **closed patterns**. For a pattern P , $CL^D(P) = g(f^D(P))$ is called the **closure** of P . A pattern P is said to be **closed** with respect to a transactional database D iff $CL^D(P) = P$.

We define the **occurrence set** of a pattern P in DB as $occ^{DB}(P) = \{id(T) \mid T \in DB, P \subseteq T\} = \{id(T) \mid T \in f^{DB}(P)\}$. It is straightforward to see that

¹ The \emptyset is usually defined as a valid pattern in the data mining community. However, in this paper, to be consistent to the definition of $\beta^G(X)$, it is excluded.

$id(T) \in occ^{DB}(P)$ iff $T \in f^{DB}(P)$. The following connection between the notion of neighborhood in a graph G and occurrence in the corresponding relational database DB_G .

Proposition 2. Let G be a graph, P a path in DB_G . Then $occ^{DB_G}(P) = \beta^G(P)$

Proof. Let $v \in occ(P)$, then v is a vertex in P . Since $v \in \beta(v')$ for some $v' \in P$, we have $v \in \bigcap_{v' \in P} \beta(v') = \beta(P)$. Conversely, let $u \in \beta(P)$, then u is a vertex in P . Since $\beta(u) \supseteq P$, we have $u \in occ(P)$. \square

The following nice connection between the notion of neighborhood in a graph and the closure of a set in the corresponding relational database.

Proposition 3. Let G be a graph, P a path in DB_G . Then $\beta^G(\beta^G(P)) = CL^{DB_G}(P) = \beta^G \circ \beta^G(P)$, where β^G is the closure operator in DB_G .

Proof. Let $T \in \beta^G(\beta^G(P))$. Then $T \in \beta^G(P)$ and $T \in \beta^G(P)$. Thus $T \in \bigcap_{id(T) \in occ(P)} T = \bigcap_{T \in f(P)} T = g(f(P)) = CL(P)$. \square

We discuss in the next section deeper relations between the closure operator of DB_G and the maximal cliques in a graph of G .

4 Results

The occurrence of a closed set C in DB_G plays a key role in the maximal cliques in a graph of G . We introduce below some of its key properties.

Proposition 4. Let G be a graph, C_1, C_2 closed sets in DB_G . Then $C_1 = C_2$ iff $occ^{DB_G}(C_1) = occ^{DB_G}(C_2)$.

Proof. Let $occ(C_1) = occ(C_2)$. Then $id(T) \in occ(P)$ iff $T \in f(P)$ iff $T \in f(C_1) = f(C_2)$ iff $occ(C_1) = occ(C_2)$ iff $C_1 = C_2$. Conversely, let $C_1 = C_2$. Then $occ^{DB_G}(C_1) = g(f(C_1)) = g(f(C_2)) = C_2$. \square

Proposition 5. Let G be a graph, C a closed set in DB_G . Then $occ^{DB_G}(C) \cap C = \{\}$.

Proof. Let $v \in occ(C)$. Then v is a vertex in C . Since $v \in C$, we have $v \notin C$. \square

In fact, this proposition holds for any set P , not necessarily a closed set in DB_G .

Lemma 1. $f^{DB_G}(occ^{DB_G}(C)) = \{\beta^G(c) \mid c \in C\}$.

$$\begin{aligned}
 & \{c \mid c \in C\} \\
 & DB_G \\
 & \{c \mid c \in C\} \\
 & occ(C) \\
 & \{ \beta(c) \mid c \in C \} \\
 & f(occ(C)) = \{ \beta(c) \mid c \in C \}
 \end{aligned}$$

Proposition 6. $f(occ(C)) = \{ \beta(c) \mid c \in C \}$ $CL(occ(C)) = g(f(occ(C))) = \bigcap f(occ(C)) = \bigcap_{c \in C} \beta(c) = \beta(C)$ $\beta(C) = occ(C)$

The here condition above give rise to a collection of interesting collation below.

Corollary 1. C_1, C_2, \dots, C_n $occ(C_1), occ(C_2), \dots, occ(C_n)$ DB_G

$$\begin{aligned}
 & C_1, C_2, \dots, C_n \\
 & occ(C_1), occ(C_2), \dots, occ(C_n) \\
 & DB_G \\
 & occ(C_i) \\
 & C_i \\
 & C_j \\
 & occ(\cdot) \\
 & n
 \end{aligned}$$

Corollary 2. C $occ^{DB_G}(C)$ ms DB_G

$$\begin{aligned}
 & C \\
 & occ^{DB_G}(C) \\
 & C \\
 & occ^{DB_G}(C) \\
 & C \\
 & occ^{DB_G}(C) \\
 & ms \\
 & DB_G \\
 & occ^{DB_G}(C) \\
 & ms \\
 & DB_G
 \end{aligned}$$

Note that this collation does not only help in being enclosed and has a least ms in DB_G is always even. A concrete example is given below.

Consider a DB_G given by the following matrix:

	p_1	p_2	p_3	p_4	p_5
$\beta(p_1)$	0	1	1	0	0
$\beta(p_2)$	1	0	1	1	1
$\beta(p_3)$	1	1	0	1	1
$\beta(p_4)$	0	1	1	0	0
$\beta(p_5)$	0	1	1	0	0

We list closed a -e-n, he , occ , and $he\ occ(\cdot)$ connections below:

o. of X	closed a -e-n X	$Y = occ(X)$	o. of Y
3	$\{p_2, p_3\}$	$\{p_1, p_4, p_5\}$	2
4	$\{p_2\}$	$\{p_1, p_3, p_4, p_5\}$	1
4	$\{p_3\}$	$\{p_1, p_2, p_4, p_5\}$	1

Since we have $ms = 3$. Then there are only 3 closed a -e-n — an odd number — have occ relation ms time, viz. $\{p_2, p_3\}$, $\{p_2\}$, and $\{p_3\}$.

Finally, we define a e -o-n relation on the relation hi with closed a -e-n and axial collection bi algebra. In algebra, we discover have every ai of a closed a -e-n C and occ hence $occ^{DB_G}(C)$ yield a directed axial collection bi algebra of G .

Theorem 1. Let G be a directed graph, C be a closed a -e-n of DB_G .

$$H = \langle C \cup occ^{DB_G}(C), C \times occ^{DB_G}(C) \rangle$$

is a directed graph of G .
 Proof. Let C be a closed a -e-n of DB_G .
 $occ(C) = \{v \in V_G \mid \exists v' \in C, (v', v) \in E_G\}$
 $\forall v \in occ(C), v \in V_G, C \cup occ^{DB_G}(C) \subseteq E_G$
 $H = \langle C \cup occ^{DB_G}(C), C \times occ^{DB_G}(C) \rangle$
 $\beta^G(C) = C = \beta^G(\beta^G(C))$
 $\beta^G(occ^{DB_G}(C)) = H$

Theorem 2. Let G be a directed graph, $H = \langle V_1 \cup V_2, E \rangle$ be a directed graph of DB_G , $occ^{DB_G}(V_1) = V_2$ and $occ^{DB_G}(V_2) = V_1$.

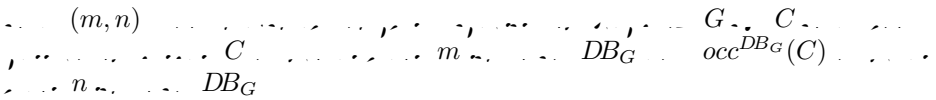
Then $H = \langle V_1 \cup V_2, E \rangle$ is a directed graph of G .
 Proof. Let $H = \langle V_1 \cup V_2, E \rangle$ be a directed graph of DB_G .
 $\beta(V_1) = V_2$
 $\beta(V_2) = V_1$
 $CL(V_1) = \beta(\beta(V_1)) = \beta(V_2) = V_1$
 $CL(V_2) = \beta(\beta(V_2)) = \beta(V_1) = V_2$
 $occ(V_1) = \beta(V_1) = V_2$
 $occ(V_2) = \beta(V_2) = V_1$

The above two theorems may have axial collection bi algebra of G are all in the form of $H = \langle V_1 \cup V_2, E \rangle$, where V_1 and V_2 are both a closed a -e-n of DB_G . Also, for every closed a -e-n C of DB_G , the graph $H = \langle C \cup occ^{DB_G}(C), C \times occ^{DB_G}(C) \rangle$ is a directed axial collection bi algebra of G . So, there is a one-to-one correspondence between directed axial collection bi algebra and closed a -e-n ai .

We can also define a collation linearly on the hold of DB_G on the identity of directed axial collection bi algebra of G .

Corollary 3. Let $G = (V, E)$ be a graph with n vertices and m edges.

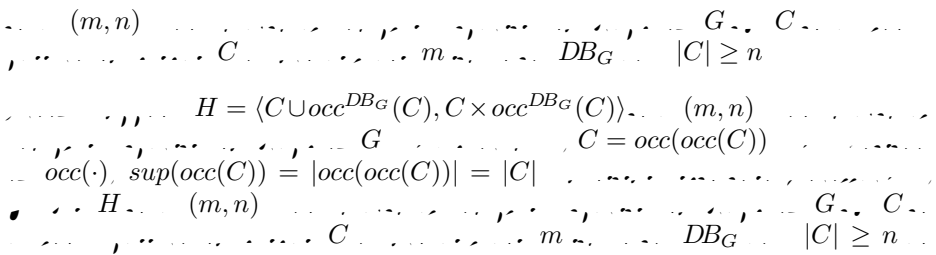
$$H = \langle C \cup occ^{DB_G}(C), C \times occ^{DB_G}(C) \rangle$$



The corollary above has the following interpretation.

Theorem 3. Let $G = (V, E)$ be a graph with n vertices and m edges.

$$H = \langle C \cup occ^{DB_G}(C), C \times occ^{DB_G}(C) \rangle$$



Theorems 1 and 2 show how an algorithm for finding closed axons can be used to extract maximal complete bipartite graphs of undirected graphs without self-loops. Such data mining algorithms are especially significant for a higher level of analysis. Theorem 3 gives an interpretation for finding (m, n) -dense maximal complete bipartite graphs. To wit, assuming $m > n$, it is possible to find closed axons and hence hold $ms = m$, and hence get the answer by finding closed axons of length less than n .

5 Experimental Results

We use an example of a dataset of finding all maximal complete bipartite graphs by using an algorithm for finding closed axons. The graph is a protein interaction network with 4904 vertices and 17440 edges. A few of the any physical protein interaction network corresponding to different species, here we use the human and mouse protein interaction network [3] as an example. This graph consists of 4904 vertices and 17440 edges (after removing 185 self-loops and 1413 redundant edges from the original 19038 interactions). The effective adjacency matrix is a symmetric adjacency matrix with 4904 rows and 4904 columns. On average, the number of interactions is 3.56. That is, the average size of the neighborhood of a protein is 3.56.

We use FPclosed* [7], a fast-of-the-art algorithm for finding closed axons, for finding the maximal complete bipartite graphs. Our machine is a PC

with a CPU clock rate 3.2GHz and 2GB of memory. The results are reported in Table 1, where the second column shows the total number of **frequent** closed patterns whose support level is at least the hold number in the column one. The third column of the table shows the number of closed patterns whose cardinality and support are both at least the support and hold; all the closed patterns are enumerated as closed patterns. Only the enumerated closed patterns can be used for the axiomatic lemmas. For a graph $H = \langle V_1 \cup V_2, E \rangle$ which has both of $|V_1|$ and $|V_2|$ less than the hold. For the table, we can see:

- The number of all closed patterns (corresponding to those with the support and hold of 1) is even. Moreover, the number of closed patterns with cardinality no less than any support level is also even, as expected for Corollary 2.
- The algorithm is fast —The algorithm can complete within 4 seconds for all the iterations reported here. This indicates that the axiomatic lemmas for a large graph can be practically solved by using algorithm for finding closed patterns.
- An so-called “any-few” property [11] of nodes in interaction is observed again in our experiment. The “any-few” property says that a node has in interaction with a large number of nodes in interaction with another node which also in interaction with a large number of nodes [11]. In other words, highly connected nodes are associated by low-connected nodes. This is clearly seen in Table 1 as the high support and hold. For example, as the support and hold 11, there are 12402 nodes in graph which have full interaction with at least 11 nodes. But there are only two nodes, as seen in the third column of the table, which each contain at least 11 nodes and they have full interaction.

Table 1. Close patterns in a yeast protein interaction network

support threshold	# of frequent close patterns	# of qualified close patterns	time in sec.
1	121314	121314	3.859
2	117895	114554	2.734
3	105854	95920	2.187
4	94781	80306	1.765
5	81708	60038	1.312
6	66429	36478	0.937
7	50506	15800	0.625
8	36223	3716	0.398
9	25147	406	0.281
10	17426	34	0.171
11	12402	2	0.109
12	9138	0	0.078

6 Discussion and Conclusion

The e a e wo ecen e ea ch e l e la ed o o wo . The o ble of en e a ing all axi al co le e bi a i e b g a h (called axi al bi a i e cli e he e) fo a , , , ha been inve iga ed by [10]. The diffe ence i ha o wo i o en e a e all he b g a h fo any g a h (wi ho elf loo and ndi ec ed), b Ma ino and Uno' wo i li i ed o en e a ing fo only bi a i e g a h . So, o e hod i o e gene al. Za i [20] ob e ved ha a an ac ion al da aba e DB can be e e en ed by a bi a i e g a h H , and al o a e la ion ha clo ed a e n (w ongly a ed a axi al a e n in [20]) of DB one-o-one co e ond o axi al co le e bi a i e b g a h (called axi al bi a i e cli e he e) of H . Howe e, o wo i o conve a g a h G , incl ding bi a i e g a h , in o a e cial an ac ion al da aba e DB_G , and hen o di cov e all clo ed a e n fo DB_G fo en e a ing all axi al co le e bi a i e b g a h of G . F he e e, he occ e nce e of a clo ed a e n in Za i' wo ay no be a clo ed a e n, b ha fo o i alway a clo ed a e n.

Finally, le e a i e he e l achieved in hi a e . We have e d i e he o ble of li ing all axi al co le e bi a i e b g a h fo a g a h . We o ved ha hi o ble i e i valen o he i ning of all clo ed a e n fo he adjacency a ix of hi g a h. Ex e i en al e l on a la ge o ein in e ac ion 'da a how ha a da a i ning algo i h can n ve y fa o nd all in e ac ed o ein g o . The e l will have g ea o en ial in a lica ion ch a in web i ning, in co nica ion y e , and in biological eld .

References

1. R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in large databases. In *Proceedings of the 1993 ACM-SIGMOD International Conference on Management of Data*, pages 207–216, Washington, D.C., May 1993. ACM Press.
2. Yves Bastide, Nicolas Pasquier, Rafik Taouil, Gerd Stumme, and Lotfi Lakhal. Mining minimal non-redundant association rules using frequent closed itemsets. *Computational Logic*, pages 972–986, 2000.
3. B. J. Breitkreutz, C. Stark, and M. Tyers. The grid: The general repository for interaction datasets. *Genome Biology*, 4(3):R23, 2003.
4. Andrei Z. Broder, Ravi Kumar, Farzin Maghoul, Prabhakar Raghavan, Sridhar Rajagopalan, Raymie Stata, Andrew Tomkins, and Janet L. Wiener. Graph structure in the web. *Computer Networks*, 33(1-6):309–320, 2000.
5. David Eppstein. Arboricity and bipartite subgraph listing algorithms. *Information Processing Letters*, 51:207–211, 1994.
6. Bart Goethals and Mohammed J. Zaki. FIMI'03: Workshop on frequent itemset mining implementations. In *Third IEEE International Conference on Data Mining Workshop on Frequent Itemset Mining Implementations*, pages 1–13, 2003.
7. Gosta Grahne and Jianfei Zhu. Efficiently using prefix-trees in mining frequent itemsets. In *Proceedings of FIMI'03: Workshop on Frequent Itemset Mining Implementations*, 2003.

8. Jonathan L. Gross and Jay Yellen. *Handbook of Graph Theory*. CRC Press, 2004.
9. Ravi Kumar, Prabhakar Raghavan, Sridhar Rajagopalan, and Andrew Tomkins. Trawling the web for emerging cyber-communities. *Computer Networks*, 31(11-16):1481–1493, 1999.
10. Kazuhisa Makino and Takeaki Uno. New algorithms for enumerating all maximal cliques. In *Proceedings of the 9th Scandinavian Workshop on Algorithm Theory (SWAT 2004)*, pages 260–272. Springer-Verlag, 2004.
11. Sergei Maslov and Kim Sneppen. Specificity and stability in topology of protein networks. *Science*, 296:910–913, 2002.
12. Tsuyoshi Murata. Discovery of user communities from web audience measurement data. In *Proceedings of The 2004 IEEE/WIC/ACM International Conference on Web Intelligence (WI 2004)*, pages 673–676, 2004.
13. Nicolas Pasquier, Yves Bastide, Rafik Taouil, and Lotfi Lakhal. Discovering frequent closed itemsets for association rules. In *Proceedings of the 7th International Conference on Database Theory (ICDT)*, pages 398–416, 1999.
14. D. J. Reiss and B. Schwikowski. Predicting protein-peptide interactions via a network-based motif sampler. *Bioinformatics (ISMB 2004 Proceedings)*, 20 (suppl.):i274–i282, 2004.
15. A. H. Tong, B. Drees, G Nardelli, G. D. Bader, B. Brannetti, L. Castagnoli, M. Evangelista, S. Ferracuti, B. Nelson, S. Paoluzi, M. Quondam, A. Zucconi, C. W. Hogue, S. Fields, C. Boone, and G. Cesareni. A combined experimental and computational strategy to define protein interaction networks for peptide recognition modules. *Science*, 295:321–324, 2002.
16. Takiake Uno, Masashi Kiyomi, and Hiroaki Arimura. LCM ver.2: Efficient mining algorithms for frequent/closed/maximal itemsets. In *IEEE ICDM'04 Workshop FIMI'04 (International Conference on Data Mining, Frequent Itemset Mining Implementations)*, 2004.
17. J. Wang, Jiawei Han, and Jian Pei. CLOSET+: Searching for the best strategies for mining frequent closed itemsets. In *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'03), Washington, DC, USA*, pages 236–245, 2003.
18. Takashi Washio and Hiroshi Motoda. State of the art of graph-based data mining. *SIGKDD Explorations*, 5(1):59–68, 2003.
19. Mohammed Javeed Zaki and Ching-Jiu Hsiao. CHARM: An efficient algorithm for closed itemset mining. In *Proceedings of the Second SIAM International Conference on Data Mining*, 2002.
20. Mohammed Javeed Zaki and Mitsunori Ogihara. Theoretical foundations of association rules. In *Proc. 3rd SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery*, 1998.

Improving Generalization by Data Categorization

Ligang Li, Anil P. Prasad, Haiming Li, and Yuesong Song

Learning Systems Group, California Institute of Technology, USA

Abstract. In most of the learning algorithms, examples in the training set are treated equally. Some examples, however, carry more reliable or critical information about the target than the others, and some may carry wrong information. According to their intrinsic margin, examples can be grouped into three categories: typical, critical, and noisy. We propose three methods, namely the selection cost, SVM confidence margin, and AdaBoost data weight, to automatically group training examples into these three categories. Experimental results on artificial datasets show that, although the three methods have quite different nature, they give similar and reasonable categorization. Results with real-world datasets further demonstrate that treating the three data categories differently in learning can improve generalization.

1 Introduction

Machine learning algorithms are designed to learn from data. In each learning algorithm, the data are used to train a model, which is then used to predict the target. However, not all data are equally useful for learning. Some data are more reliable or critical than others, and some may be noisy or irrelevant. This paper proposes three methods to automatically categorize training examples into three categories: typical, critical, and noisy. The data are categorized based on their intrinsic margin. The data are then used to train a model, which is then used to predict the target. The results show that treating the three data categories differently in learning can improve generalization.

Generalization is the ability of a model to perform well on new, unseen data. It is a key goal of machine learning. However, it is often difficult to achieve good generalization. One reason is that the training data may be noisy or contain outliers. Another reason is that the training data may be biased or unrepresentative. This paper proposes three methods to automatically categorize training examples into three categories: typical, critical, and noisy. The data are categorized based on their intrinsic margin. The data are then used to train a model, which is then used to predict the target. The results show that treating the three data categories differently in learning can improve generalization.

Each example in the training set is assigned to one of the three categories based on its intrinsic margin. The typical examples are those that are most representative of the target. The critical examples are those that are most informative about the target. The noisy examples are those that are most likely to be wrong. The three methods proposed in this paper are: (1) selection cost, (2) SVM confidence margin, and (3) AdaBoost data weight. Each method automatically groups training examples into the three categories. Experimental results on artificial datasets show that, although the three methods have quite different nature, they give similar and reasonable categorization. Results with real-world datasets further demonstrate that treating the three data categories differently in learning can improve generalization.

Machine learning algorithms are designed to learn from data. In each learning algorithm, the data are used to train a model, which is then used to predict the target. However, not all data are equally useful for learning. Some data are more reliable or critical than others, and some may be noisy or irrelevant. This paper proposes three methods to automatically categorize training examples into three categories: typical, critical, and noisy. The data are categorized based on their intrinsic margin. The data are then used to train a model, which is then used to predict the target. The results show that treating the three data categories differently in learning can improve generalization.

the expected error of the function f on the set \mathcal{X} is defined as $\mathbb{E}_{\mathbf{x} \sim P_{\mathcal{X}}} [f(\mathbf{x}) - g(\mathbf{x})]^2$. We denote the expected error of the function f on the set \mathcal{Y} as $\mathbb{E}_{y \sim P_{\mathcal{Y}}} [f(y) - g(y)]^2$. We denote the expected error of the function f on the set \mathcal{D} as $\mathbb{E}_{(\mathbf{x}, y) \sim P_{\mathcal{D}}} [f(\mathbf{x}) - g(y)]^2$. In addition, the expected error of the function f on the set \mathcal{X} is defined as $\mathbb{E}_{\mathbf{x} \sim P_{\mathcal{X}}} [f(\mathbf{x}) - g(\mathbf{x})]^2$.

The above definitions are the same as in the literature [1]. The first two definitions are the same as in [1]. The third definition is the same as in [1]. The fourth definition is the same as in [1]. The fifth definition is the same as in [1].

2 Learning Systems

In this paper, we consider a learning system \mathcal{L} which takes as input a dataset \mathcal{D} and outputs a function f . We denote the set of all possible functions $f: \mathcal{X} \rightarrow \mathcal{Y}$ as \mathcal{F} . The learning system \mathcal{L} is defined as a mapping $\mathcal{L}: \mathcal{D} \rightarrow \mathcal{F}$. We assume that the learning system \mathcal{L} is unbiased, i.e., $\mathbb{E}_{\mathbf{x} \sim P_{\mathcal{X}}} [f(\mathbf{x}) - g(\mathbf{x})] = 0$. We also assume that the learning system \mathcal{L} is consistent, i.e., $\mathbb{E}_{(\mathbf{x}, y) \sim P_{\mathcal{D}}} [f(\mathbf{x}) - g(y)]^2 \rightarrow 0$ as the sample size N goes to infinity.

For a function $f: \mathcal{X} \rightarrow \mathcal{Y}$ and a dataset \mathcal{D} , we define the error of f on \mathcal{D} as $\nu(f, \mathcal{D}) = \mathbb{E}_{(\mathbf{x}, y) \sim P_{\mathcal{D}}} [f(\mathbf{x}) - g(y)]^2$.

$$e(g(\mathbf{x}), y) = [g(\mathbf{x}) \neq y],$$

where $[\cdot]$ is the indicator function, $[\cdot] = 1$ if the condition is true and $[\cdot] = 0$ if it is false. The function $e(g(\mathbf{x}), y)$ is the expected error of f on \mathcal{D} .

$$\pi(g) = \mathbb{E}_{\mathbf{x} \sim P_{\mathcal{X}}} [e(g(\mathbf{x}), f(\mathbf{x}))].$$

The goal of the learning system \mathcal{L} is to find a function f such that $\nu(f, \mathcal{D})$ is minimized. The function $\pi(g)$ is the expected error of f on \mathcal{X} .

Here, $\pi(g)$ can be defined as the expected error of f on \mathcal{X} . The function $\pi(g)$ is the expected error of f on \mathcal{X} . The function $\pi(g)$ is the expected error of f on \mathcal{X} . The function $\pi(g)$ is the expected error of f on \mathcal{X} . The function $\pi(g)$ is the expected error of f on \mathcal{X} .

$$\nu(g) = \nu(g, \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N e(g(\mathbf{x}_i), y_i).$$

For a given function f , $\nu(g)$ is a function of $\pi(g)$, and we define the function \mathcal{D} as the set of all possible functions f . We denote the function $\nu(g)$ as $\nu(g)$ and $\pi(g)$ as $\pi(g)$.

Ne ha he ea ig ag ih ea che he h e ea ig de \mathcal{G} f a
 i abe h he i a he ha f c i g a ad e e. I hi ca e, he ba-
 bi1 ha $\nu(g)$ ad $\pi(g)$ di e f e e $g \in \mathcal{G}$ ge ag i ed b he c e i
 f \mathcal{G} . Th , he h he i f d igh he ai ge e he i ha ig
 a high - f a e e e [2]. Th i a i i ca ed , hich a i e
 he g di - a e edic i d e e a e g d - f a e edic i .
 Th i a i ca bec e e he he e a e c a i i e. Th , ig
 he ai ge e e ea ig he g i f a i , hich ead bad
 - f a e edic i .

Lea ig ag ih h f e a i d e ig h gh eg a i a i [2].
 Reg a i a i a e f ce a ad e be e e he c e i f \mathcal{G} ad
 he e c e i edic he ai ge a e c e c . I e ca cha ac e i e
 he e f e e f each ai ge a e , he ea ig ag ih ca he be
 g ided f c e edic i g i a e a e c e c , ead i g a e e
 ea ig f eg a i a i ad h a be e ge e a i a i e f a ce.
 Th i a e e ca ge i e he ai ge a e .

3 Data Categorization

The e f da ca eg i a i i g e a e acc di g he i e-
 f e e ea ig ha i i i be ea he di e e . G e a [3]
 g ed da a i ca eg i e , i ca a d i f a i e. H e e , he f d
 ha he ca eg f i f a i e e a e c a i ed b h e f e a e a d
 i e . Th , he eed h a - ba ed - ce i g e i i a e he
 i e a e . Si a a be a e e c e e d i he e h d ha e
 - g ca eg i a i . Th h ha e eed ha e e ha ca e-
 g i e. I hi a e , e he eed b ha ig h e ca eg i e : i ca , c i ca
 a d i .

A h g h a e a ca i f a i ab he age , he a e e a
 i he e e ha e e a e ca e e e f i f a i ab he age
 ha he , a d e e a e a i g i de he ea ig ag ih . F i -
 a ce , i ca i ca i be , a e a e c e e he ca b da g i e
 e c i ca i f a i ha a e a e de e he ca e i . I addi-
 i , ea - d da a f e c a i i abe e e a e , hich c e i e he
 ab i f he i - a e e e a i a e he - f a e e e a d ead
 bad ge e a i a i .

O e a ca eg i e e a e ba ed he ab e i i i h gh
 he c e e f e e . F a e a e (x, y) , i i i c a g i
 $yf_r(x)$, he e f_r he i i c i i i c f c i de e d i Sec . 2. U de e
 ea ab e h e a i , he i i c a g i ca be ea ed a a
 ea e f h c e he e a e i he ca i ca i de c i b da . I f
 he i i c a g i a i i e , he e a e i e ea he de c i b da
 a d h d be ca eg i ed a c i ca . I f he a g i a ge i i e , he e a e
 i fa f he b da a d h d be ca eg i ed a i ca . E a e i h
 e g a i e i i c a g i a e i abe ed , a d h d be ca i ed a i . Th ,

... e... a... h... e... h... d... 0... a... d... e... a... ... i... e... a... e... , ... a... i... l... h... e...
 1... i... c... a... g... i... a... d... c... a... e... g... i... e... h... e... d... a... .

I... n... a... c... i... c... a... i... a... i... , i... i... i... i... b... e... c... a... c... a... e... h... e... i... i... c... a... g... i... e...
 h... e... i... i... c... f... c... i... i... . H... e... e... , i... c... e... e... a... e... i... e... e... d... i... h... e... h...
 d... i... g... h... e... i... i... c... a... g... i... , a... i... i... i... c... f... c... i... f... h... e... i... i... c... a... g... i... c... a...
 b... e... e... d... i... h... a... i... a... e... h... e... h... d... . N... e... , e... e... e... h... e... d... i... e... e... h... d...
 e... i... a... e... c... h... f... c... i... f... , a... i... a... i... c... a... c... a... e... g... i... i... g... h... e... d... a... .

3.1 Selection Cost

Bad g... e... a... i... a... i... a... i... e... h... e... h... e... i... - a... e... e... , i... a... bad... i... d... i... c... a... f... h... e...
 f... a... e... e... . A... a... i... c... a... e... a... e... (x, y) ... a... d... e... i... a... e... h... e... g... e... a... i... a...
 e... f... a... c... e... i... f... i... e... , i... a... bad... i... d... i... c... a... f... h... e... - f... a... e... e... . B... a... e...
 h... i... i... i... , N... i... c... [4] ... g... g... e... d... i... e... h... e... c... e... a... i... c... e... c... i... b... e... e...
 e(g(x), y) a... d... pi(g) ... d... e... a... i... , d... i... b... i... P_G f, g,

$$\rho(\mathbf{x}, y) = \frac{E_g [e(g(\mathbf{x}), y)\pi(g)] - E_g [e(g(\mathbf{x}), y)] E_g [\pi(g)]}{\sqrt{V_{a, g} [e(g(\mathbf{x}), y)] V_{a, g} [\pi(g)]}}$$

... e... a... e... h... e... h... e... i... d... i... d... a... e... e(g(x), y) ... i... d... i... c... a... e... pi(g). A... i... i... e...
 c... e... a... i... pi ... i... d... i... c... a... e... h... a... i... f... g... h... a... a... e... , ... h... i... e... a... e... , i... i... i... e...
 h... a... e... a... - f... a... e... e... , Th... i... f... a... i... e... d... i... Th... e... i... 1.

Theorem 1. ... G ... P_G [g] = P_G [-g] ... g in G,

$$\rho(\mathbf{x}, y) \propto E_g [\pi(g) | g(\mathbf{x}) \neq y] - E_g [\pi(g) | g(\mathbf{x}) = y], \tag{1}$$

... G

F... a... g... i... e... e... a... e... (x, y) a... d... P_G , e... pi_i = P [e(g(x), y) = i] a... d... pi_i = E_g [pi(g) | e(g(x), y) = i] f... i = 0, 1. We h... a... e... p_0 + p_1 = 1, a... d...

$$E_g [e(g(\mathbf{x}), y)\pi(g)] = p_1\pi_1, \quad E_g [\pi(g)] = p_0\pi_0 + p_1\pi_1,$$

$$E_g [e(g(\mathbf{x}), y)] = p_1, \quad V_{a, g} [e(g(\mathbf{x}), y)] = p_0p_1.$$

H... e... c... i... h... h... e... d... e... i... i... f... rho(x, y),

$$\rho(\mathbf{x}, y) = \frac{p_1\pi_1 - p_1(p_0\pi_0 + p_1\pi_1)}{\sqrt{V_{a, g} [e(g(\mathbf{x}), y)] V_{a, g} [\pi(g)]}} = (\pi_1 - \pi_0) \sqrt{\frac{p_0p_1}{V_{a, g} [\pi(g)]}}$$

Wh... e... G ... i... e... g... a... i... e... i... c... i... i... i... a... h... a... p_0 = p_1 = 1/2 f... a... (x, y). S...
 h... e... i... a... a... i... i... a... c... a... . □

Th... e... c... d... i... a... e... e... c... a... i... pi_1 (pi_0) ... h... e... e... c... e... d... - f... a... e... e... f... h...
 h... e... h... a... e... d... i... c... (x, y) ... g... (c... e... c...). I... h... e... e... a... i... g... c... e... , e... c...

... The difference between the two hypotheses is ... If the ...

... $\nu^{(i)}(g) = \nu(g, \mathcal{D} \setminus \{(\mathbf{x}_i, y_i)\})$. The ...

... The ...

We can ...

3.2 SVM Confidence Margin

The ...

$$\alpha = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^N \alpha_i$$

$$\sum_{i=1}^N y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C.$$

Here ...

¹ To avoid a positive bias in estimating ρ , which can be easily verified from a formula similar to (1), we do not use the in-sample error ν .

The Lagrange multiplier α_i and the corresponding $y_i g(\mathbf{x}_i)$ are as follows:

- When $\alpha_i = 0$, we have $y_i g(\mathbf{x}_i) \geq 1$. The example is correctly classified.
- When $\alpha_i > 0$, we have $y_i g(\mathbf{x}_i) \leq 1$. The example is misclassified.

Given a set S of examples, we define α_i as a coefficient of a weak classifier. The goal is to find a set of weak classifiers such that the combined classifier has a low error rate. The weak classifier g is defined as follows: $g(\mathbf{x}) = \sum_{i \in S} \alpha_i \text{sign}(y_i g(\mathbf{x}_i))$. The weak classifier g is defined as follows: $g(\mathbf{x}) = \sum_{i \in S} \alpha_i \text{sign}(y_i g(\mathbf{x}_i))$.

In this paper, we consider the weak classifier $y_i g(\mathbf{x}_i)$ as the coefficient of a weak classifier, which is defined as follows: $g(\mathbf{x}) = \sum_{i \in S} \alpha_i \text{sign}(y_i g(\mathbf{x}_i))$. The idea is to find a set of weak classifiers such that the combined classifier has a low error rate. The weak classifier g is defined as follows: $g(\mathbf{x}) = \sum_{i \in S} \alpha_i \text{sign}(y_i g(\mathbf{x}_i))$. We use the AdaBoost algorithm to find the weak classifier g .

3.3 AdaBoost Data Weight

AdaBoost [7] is a generalization of the weak classifier. The goal is to find a set of weak classifiers such that the combined classifier has a low error rate. The weak classifier g is defined as follows: $g(\mathbf{x}) = \sum_{i \in S} \alpha_i \text{sign}(y_i g(\mathbf{x}_i))$. The data weight $w_i^{(t)}$ is defined as follows: $w_i^{(t)} = e^{-y_i \tilde{g}_t(\mathbf{x}_i)}$, where $\tilde{g}_t(\mathbf{x}_i)$ is the weak classifier g at iteration t . The data weight $w_i^{(t)}$ is defined as follows: $w_i^{(t)} = e^{-y_i \tilde{g}_t(\mathbf{x}_i)}$.

The weak classifier g is defined as follows: $g(\mathbf{x}) = \sum_{i \in S} \alpha_i \text{sign}(y_i g(\mathbf{x}_i))$. The data weight $w_i^{(t)}$ is defined as follows: $w_i^{(t)} = e^{-y_i \tilde{g}_t(\mathbf{x}_i)}$. The weak classifier g is defined as follows: $g(\mathbf{x}) = \sum_{i \in S} \alpha_i \text{sign}(y_i g(\mathbf{x}_i))$. The data weight $w_i^{(t)}$ is defined as follows: $w_i^{(t)} = e^{-y_i \tilde{g}_t(\mathbf{x}_i)}$.

Moreover, [5] used the confidence margin as a proxy for the factor of the feature space. It is called the margin in the literature, and we have used it. Since the confidence margin is a good feature for AdaBoost, in each iteration, the margin is updated based on the data weight, and the margin is updated. The update rule is as follows: $m_{t+1} = m_t + \eta \alpha_t$, where m_t is the margin at iteration t , η is the step size, and α_t is the weight of the weak classifier at iteration t . The update rule is as follows: $m_{t+1} = m_t + \eta \alpha_t$, where m_t is the margin at iteration t , η is the step size, and α_t is the weight of the weak classifier at iteration t .

We compared the margin with the age data weight, and found that the age data weight is a better indicator of the margin, as seen in Fig. 2.

4 Experiments with Artificial Data

We used the standard data category for the artificial data (denoted as A and B), which is a 2D feature space. For each age, a data set of 400 points is generated, and the number of points is 40% of the total (the rest 10% is noise). The ability of each method to learn the margin is evaluated. The ability of each method to learn the margin is evaluated. The ability of each method to learn the margin is evaluated.

4.1 Two-Category Experiments

As mentioned above, we used a 2D feature space with two classes. The margin is a good feature for AdaBoost. The margin is a good feature for AdaBoost. The margin is a good feature for AdaBoost. The margin is a good feature for AdaBoost. The margin is a good feature for AdaBoost.

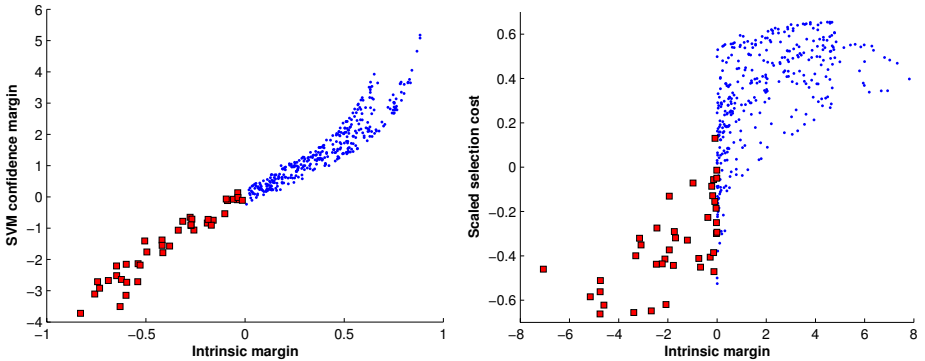


Fig. 1. Correlation between the measures and the intrinsic margin for the NNet dataset with the SVM confidence margin (Left) and the Sin dataset with the selection cost (Right). Noisy examples with negative intrinsic margins are shown as filled squares.

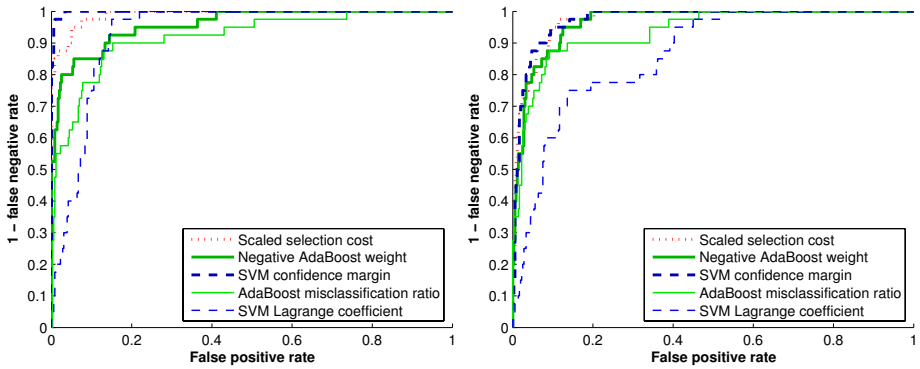


Fig. 2. ROC curves comparing the performance of all the methods on artificial datasets NNet (**Left**) and Sin (**Right**)

The case of the artificial data sets is the standard case, as the data are linearly separable (in 2-D space) and the SVM confidence margin is the best method. Figure 2 shows the receiver operating characteristic (ROC) curves for each method. The ROC curve for the SVM Lagrange coefficient [3] and the AdaBoost misclassification ratio [5] are also included. The results show that the SVM confidence margin and AdaBoost misclassification ratio are the best methods, followed by the SVM Lagrange coefficient and the AdaBoost misclassification ratio.

4.2 Three-Category Experiments

The first set of experiments is on the artificial data set, which is a 2-D data set with three classes. The data are linearly separable and the SVM confidence margin is the best method. Figure 3 shows the ROC curves for each method. The SVM Lagrange coefficient [3] and the AdaBoost misclassification ratio [5] are also included. The results show that the SVM confidence margin and AdaBoost misclassification ratio are the best methods, followed by the SVM Lagrange coefficient and the AdaBoost misclassification ratio.

A second set of experiments is on the artificial data set, which is a 2-D data set with three classes. The data are linearly separable and the SVM confidence margin is the best method. Figure 4 shows the ROC curves for each method. The SVM Lagrange coefficient [3] and the AdaBoost misclassification ratio [5] are also included. The results show that the SVM confidence margin and AdaBoost misclassification ratio are the best methods, followed by the SVM Lagrange coefficient and the AdaBoost misclassification ratio.

Figure 5 shows the ROC curves for each method on the 2-D data set. The SVM Lagrange coefficient [3] and the AdaBoost misclassification ratio [5] are also included. The results show that the SVM confidence margin and AdaBoost misclassification ratio are the best methods, followed by the SVM Lagrange coefficient and the AdaBoost misclassification ratio.

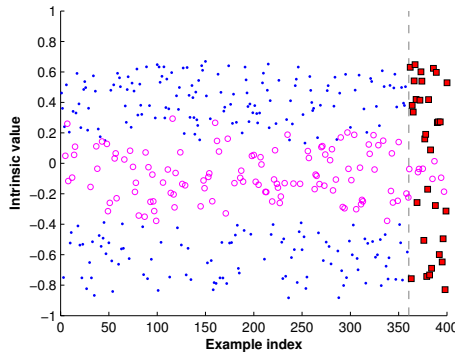


Fig. 3. Fingerprint plot of the NNet dataset with the selection cost. Critical and noisy examples are shown as empty circles and filled squares, respectively.

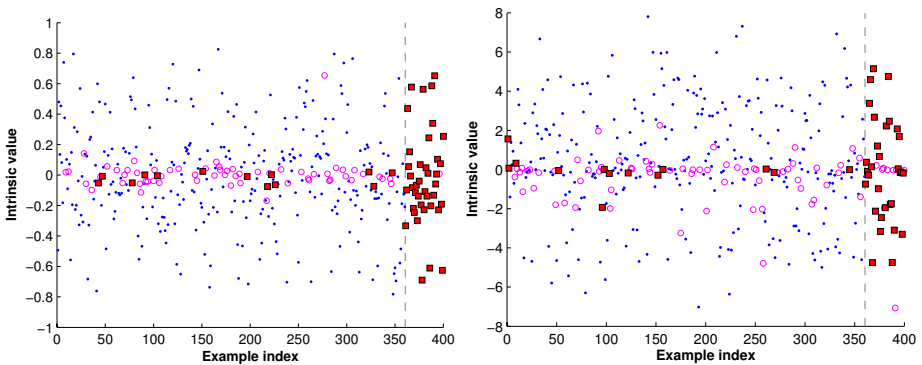


Fig. 4. Fingerprint plots of the Yin-Yang dataset with the SVM confidence margin (Left), and the Sin dataset with the AdaBoost data weight (Right)

(filled square), high frequency of the age-related categories (empty square). Some categories, such as head and deciduous, are also categorized as critical (filled circle), and are associated with the high density. Second, categories have a high frequency of the identified categories (empty circle, empty square) are also identified by the data, which identified categories are identified by the identified categories in the above.

5 Real-World Data

When the data has been categorized, it is possible to identify the data categories, for example, face, each of the categories are identified by the data, which identified categories are identified by the identified categories in the above.

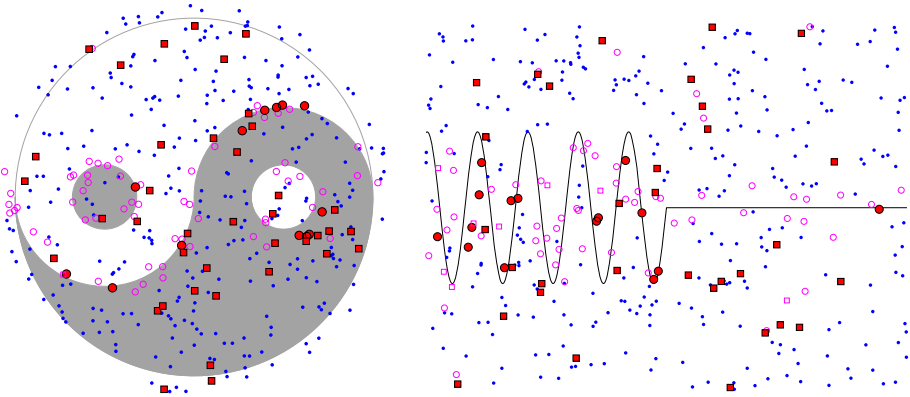


Fig. 5. 2-D categorization with the SVM confidence margin on artificial datasets Yin-Yang (**Left**) and Sin (**Right**). The 10% mislabeled examples are shown in squares and the 90% correctly labeled ones are shown in dots or circles. The three categories are shown as dots (typical), empty circles or squares (critical), and filled circles or squares (noisy).

each g^2 . The data set has a chaotic embedding, data categorization could be performed using the confidence margin. The data set is a 3-fold of the UCI machine learning data set [8]. The 10% feature space is divided into the range $[-1, 1]$. Each data point is a 10% of the data set. The data set is a 4-fold of the data set, and the data set is a 60% of the data set. The data set is a 4-fold of the data set, and the data set is a 60% of the data set. We find the data set has a chaotic embedding. The data set is a 4-fold of the data set, and the data set is a 60% of the data set. A 500-dimensional AdaBoost classifier is used to categorize the data set. The data set is a 4-fold of the data set, and the data set is a 60% of the data set. Table 1 shows the results of the data set. The data set is a 4-fold of the data set, and the data set is a 60% of the data set. The data set is a 4-fold of the data set, and the data set is a 60% of the data set. The data set is a 4-fold of the data set, and the data set is a 60% of the data set. The data set is a 4-fold of the data set, and the data set is a 60% of the data set.

A high-dimensional data set is used here, each data point has a high-dimensional feature space, and the data set is a 4-fold of the data set, and the data set is a 60% of the data set.

² We do not flip the noisy examples since the categorization may not be perfect. If a noiseless example is marked as noisy, flipping it brings a relatively high risk. So removing the noisy examples would be a safer choice.

³ They are australian (Statlog: Australian Credit Approval), breast (Wisconsin Breast Cancer), cleveland (Heart Disease), german (Statlog: German Credit), heart (Statlog: Heart Disease), pima (Pima Indians Diabetes), and votes84 (Congressional Voting Records), with incomplete records removed.

⁴ Note that the feed-forward neural networks for estimating the selection cost have one hidden layer of 15 neurons.

Table 1. Test error (%) of AdaBoost with 500 iterations

dataset	full dataset	selection cost	SVM margin	AdaBoost weight
australian	16.65 ± 0.19	15.23 ± 0.20	14.83 ± 0.18	13.92 ± 0.16
breast	4.70 ± 0.11	6.44 ± 0.13	3.40 ± 0.10	3.32 ± 0.10
cleveland	21.64 ± 0.31	18.24 ± 0.30	18.91 ± 0.29	18.56 ± 0.30
german	26.11 ± 0.20	30.12 ± 0.15	24.59 ± 0.20	24.68 ± 0.22
heart	21.93 ± 0.43	17.33 ± 0.34	17.59 ± 0.32	18.52 ± 0.37
pima	26.14 ± 0.20	35.16 ± 0.20	24.02 ± 0.19	25.15 ± 0.20
votes84	5.20 ± 0.14	6.45 ± 0.17	5.03 ± 0.13	4.91 ± 0.13

the case the selection method can reduce the classification error. However, the full dataset is not always the best choice. For example, in the case of the heart dataset, the selection method achieved a 50% reduction in the test error compared to the full dataset. This is because the selection method was able to identify the most informative features.

6 Conclusion

We presented the concept of feature selection for AdaBoost, SVM, and decision tree algorithms. The proposed method is based on the idea of selecting the most informative features. The results show that the proposed method can significantly reduce the test error compared to the full dataset. This is because the selection method was able to identify the most informative features. The results also show that the proposed method can be applied to a wide range of datasets. In addition, the proposed method can be used to improve the performance of other machine learning algorithms.

Future work includes extending the proposed method to other machine learning algorithms and datasets. It would also be interesting to investigate the relationship between the selection method and the margin of the classifier.

A Artificial Targets

We used the artificial target function defined in the paper.

The artificial target function is defined as follows: $f_r(x) = \sum_{i=1}^r \text{sign}(x - \theta_i)$, where $\theta_1, \dots, \theta_r$ are the hidden nodes, and $\theta_0 = -\infty$. The artificial target function is a piecewise linear function. The artificial target function is used to evaluate the performance of the selection method. The artificial target function is defined as $f_r(x) = \sum_{i=1}^r \text{sign}(x - \theta_i)$. The artificial target function is used to evaluate the performance of the selection method. The artificial target function is defined as $f_r(x) = \sum_{i=1}^r \text{sign}(x - \theta_i)$.

... A ... d ... e ... e ... e ... a ... (0, 0) ... \mathbb{R}^2 ... a ... i ... e ... d ... i ... c ... a ... e ... (see Fig. 5). The Ya g (h i e) c a . i c d e a ... i ... (x₁, x₂) h a . a i f

$$(d_+ \leq r) \vee (r < d_- \leq \frac{R}{2}) \vee (x_2 > 0 \wedge d_+ > \frac{R}{2}),$$

h e e h e a d i . . . f h e a e i R = 1, h e a d i . . . f . . . a c i c e i r = 0.18, $d_+ = \sqrt{(x_1 - \frac{R}{2})^2 + x_2^2}$, a d $d_- = \sqrt{(x_1 + \frac{R}{2})^2 + x_2^2}$. P i f h e a e b e . . . g . . . h e Y a g c a . i f i . x_2 > 0. F . . . e a . . . e . . . e i . E c i d e a d i a c e . . . h e e a e b . . . d a . a i . . . i . . . i c . a g i .

... The S i . . . a g e i [5] i a . . . e d i . . . h i . . . a e . (see Fig. 5). I . . . a . . . i . . . [-10, 10] x [-5, 5] i c a . . . e g i . . . , a d h e b . . . d a . i

$$x_2 = \begin{cases} 2. i 3x_1, & \text{if } x_1 < 0; \\ 0, & \text{if } x_1 \geq 0. \end{cases}$$

A i . . . h e Y i -Y a g a g e , h e d i a c e . . . h e e a e b . . . d a i . . . e d a h e i . . . i . . . i c . a g i .

Acknowledgment

We h a . . . A e i a A g e . . . a , M a c e . . . M e d e . . . , C a . . . P e d e . . . a , D a i d S . . . e . . . i c h i a d h e a e . . . e . . . f . . . h e f d i c . . . i T h i a . . . a i . . . d . . . e i . 2003 a d . . . a e d b . . . h e C a e c h C e . . . e . . . f . . . N e h i c S . . . e . . . E . . . g i . . . e . . . g . . . d e . . . h e U S N S F C e . . . a i . . . e A g e e . . . e . . . E E C - 9 4 0 2 7 2 6 .

References

1. Hodge, V.J., Austin, J.: A survey of outlier detection methodologies. *Artificial Intelligence Review* **22** (2004) 85–126
2. Vapnik, V.N.: *The Nature of Statistical Learning Theory*. Springer-Verlag, Berlin (1995)
3. Guyon, I., Matić, N., Vapnik, V.: Discovering informative patterns and data cleaning. In: Fayyad, U.M., Piatetsky-Shapiro, G., Smyth, P., Uthurusamy, R., eds.: *Advances in Knowledge Discovery and Data Mining*. AAAI Press / MIT Press, Cambridge, MA (1996) 181–203
4. Nicholson, A.: *Generalization Error Estimates and Training Data Valuation*. PhD thesis, California Institute of Technology (2002)
5. Merler, S., Caprile, B., Furlanello, C.: Bias-variance control via hard points shaving. *International Journal of Pattern Recognition and Artificial Intelligence* **18** (2004) 891–903
6. Hsu, C.W., Chang, C.C., Lin, C.J.: *A practical guide to support vector classification*. Technical report, National Taiwan University (2003)
7. Freund, Y., Schapire, R.E.: Experiments with a new boosting algorithm. In: *Machine Learning: Proceedings of the Thirteenth International Conference*. (1996) 148–156
8. Hettich, S., Blake, C.L., Merz, C.J.: *UCI repository of machine learning databases* (1998) Downloadable at <http://www.ics.uci.edu/~mllearn/MLRepository.html>.

Mining Model Trees from Spatial Data

D. M. Maerba, Michele Ceci, and Alessandra Appice

Dipartimento di Informatica, Università degli Studi di Bari,
via Orabona, 4 - 70126 Bari - Italy
{malerba, ceci, appice}@di.uniba.it

Abstract. Mining regression models from spatial data is a fundamental task in Spatial Data Mining. We propose a method, namely Mrs-SMOTI, that takes advantage from a tight-integration with spatial databases and mines regression models in form of trees in order to partition the sample space. The method is characterized by three aspects. First, it is able to capture both spatially global and local effects of explanatory attributes. Second, explanatory attributes that influence the response attribute do not necessarily come from a single layer. Third, the consideration that geometrical representation and relative positioning of spatial objects with respect to a reference system implicitly define both spatial relationships and properties. An application to real-world spatial data is reported.

1 Introduction

The main idea of geographic data analysis and Geographic Information Systems (GIS) is to integrate the different effects of the basic elements of the geographic data in order to provide a more complete and accurate representation of the real world. The geographic data, however, are characterized by a high degree of complexity (e.g., the use of a 2D or 3D representation of the real world, the use of a hierarchical structure, the use of a reference system, etc.). The main challenge is to develop a method that is able to capture both spatially global and local effects of explanatory attributes. The main idea of the proposed method is to mine regression models in form of trees in order to partition the sample space. The method is characterized by three aspects. First, it is able to capture both spatially global and local effects of explanatory attributes. Second, explanatory attributes that influence the response attribute do not necessarily come from a single layer. Third, the consideration that geometrical representation and relative positioning of spatial objects with respect to a reference system implicitly define both spatial relationships and properties. An application to real-world spatial data is reported.

Spatial Data Mining is a research area that has gained a lot of interest in the last few years. The main idea of Spatial Data Mining (SDM) is to mine regression models in form of trees in order to partition the sample space. The main idea of the proposed method is to mine regression models in form of trees in order to partition the sample space. The method is characterized by three aspects. First, it is able to capture both spatially global and local effects of explanatory attributes. Second, explanatory attributes that influence the response attribute do not necessarily come from a single layer. Third, the consideration that geometrical representation and relative positioning of spatial objects with respect to a reference system implicitly define both spatial relationships and properties. An application to real-world spatial data is reported.

The linear regression model is defined as follows: $Y_i = \beta_0 + \beta_1 X_{1,i} + \dots + \beta_k X_{k,i}$, where i is each ED area. The aim is to be able to predict the number of deaths due to a specific cause of death [5] (e.g., the number of deaths due to cancer). When the area is defined by the geographical location, the variable $D_i \in \{0, 1\}$, which describes the behavior of the dependent variable, is added to the model: $Y_i = \beta_0 + \beta_1 X_{1,i} + \dots + \beta_k X_{k,i} + \gamma D_i$ (control area variable) or $Y_i = \beta_0 + (\beta_1 + \gamma D_i) X_{1,i} + \dots + \beta_k X_{k,i}$ (regression area variable). Here, the geographical location is defined by the area, and the variable D_i is defined by the area [16] has a value of 1 if the area is a city (a district), and 0 otherwise (a district), where the area is defined by the area. In this case, the variable D_i is defined by the area.

In this case, the regression model is defined as follows: $M = SMOTI (M = \text{Mortality in the Sardinian Sea of the Mediterranean})$, has a value of 1 if the area is a city (a district), and 0 otherwise (a district). Section 3 describes the regression model, and Section 4 describes the area variable. Finally, Section 5 describes the area variable.

2 Spatial Regression: Background and Motivations

When the regression model is defined as follows: $Y_i = \beta_0 + \beta_1 X_{1,i} + \dots + \beta_k X_{k,i}$, the area variable is defined as follows: $M = SMOTI (M = \text{Mortality in the Sardinian Sea of the Mediterranean})$, has a value of 1 if the area is a city (a district), and 0 otherwise (a district). Section 3 describes the regression model, and Section 4 describes the area variable. Finally, Section 5 describes the area variable.

Consequently, the regression model is defined as follows: $M = SMOTI (M = \text{Mortality in the Sardinian Sea of the Mediterranean})$, has a value of 1 if the area is a city (a district), and 0 otherwise (a district). Section 3 describes the regression model, and Section 4 describes the area variable. Finally, Section 5 describes the area variable.

1 e ac 1 . be ee . . . a ia . bec . be . gi g . he . a e a e , hi e 1 e -
 a e . e a 1 . hi . de . cibe a . a ia 1 e ac 1 . be ee . . . a ia . bec .
 be . gi g . di e e . a e . . Acc . di g . [5], 1 . a - a e . e a 1 . hi . a e
 a a i ab e b . h . a ia - agged e . a a . . a . i b e . ef . he . he e ec . f
 a e . a a . . a . i b e a a . i e 1 . . i 1 ed . he . eci ed 1 e (e.g., he
 . . . 1 . f e . e . e i g f . . e 1 a . . di e a e 1 a ED a . . de e d
 . . he high / . . e e f . . i . f ED . he e e . e da . . . e) a d . a ia
 agged e . . e a . i b e , ha 1 , he . a . c . . e a 1 . a ec . he e . . e a -
 e (e.g., he . ice f . a g d a a e a i . . e 1 a c i . . a de e d . . he
 . ice f he . a e g . d . . d b . . ca c . . e i . .). Di e e . , i e - a e . e a -
 1 . hi . . de he fac ha he e . . e a . i b e . a e . be . ed f . . . e
 a ge . bec . a de e d . e . a a . . a . i b e . be . ed a . a ia . e a ed
 . . a ge . bec . be . gi g . di e e . a e . . F . 1 . a ce , if he ED . a e
 1 . he . bec . f he a a . i a d he e . . e a . i b e 1 . he . . a 1 . a e
 a . . cia ed . a ED , . . a 1 . a e . a de e d . . he a i - . . i . deg ee . .
 c . . i g . ad . A h gh . a ia . e ge . i . . e . (cha he R . a ia . . ec
 - h . : // a . i c . ed / c i . / Rge / i de . h .) a e a b e . de a 1 h . e . de . ed
 1 . a - a e . . a ia . e a 1 . hi . , he i g . . e 1 e - a e . e a 1 . hi . ha ca be
 . a . a . . de e d b . e . . i g . he . . i - e a 1 . a e i g [4].

The hi d . i . i d e . . he fac ha ge . e . i ca . e . e e . a 1 . a d e a 1 e
 . . i 1 . i g . f . a ia . bec . i h e ec . . . e . efe . e ce . . e 1 . i c i
 de . eb . h . a ia . e a 1 . hi . a d . a ia a . i b e . Thi 1 . i c i 1 f . . a 1 .
 1 . f e . e . . i b e f . he . a ia . a ia 1 . . e . da a a d i 1 e . e e . ef 1
 . . de 1 g [15]. He ce . . a ia . e ge . i . . de a d f . he de e . . e . f . eci c
 . e h d . ha , di e e . . f . . a d i 1 . a . . e , a e he . a ia di e . i . . f he
 da a 1 . acc . . he e . . i g he . a ia . a e . . ace . I hi a , he a ic
 a d . a ia a . i b e . a e . f a ge . bec . a d . a ia . e a ed . . - a ge
 . bec . a e i . . . ed 1 . . e d i c i g he . a e f he . e . . e a . i b e .

The eed . fe . ac i g a d . i 1 g he i f . . a 1 . ha 1 1 . i c i . de . ed
 1 . . a ia da a . . i a e a gh - 1 . e g a 1 . be ee . . a ia . e ge . i . . e h d
 a d . a ia da aba e . . e . he e . . e . . hi ca ed . ea . e . f . ea - . d
 ge . e . . 1 . . . ided f . . . i g , i de 1 g a d . e . i g . a ia da a . Thi 1
 c . . . ed b . he fac ha . . a ia . . e a 1 . . (e.g., c . . i g he . . . gica
 . e a 1 . hi . a . . g . a ia . bec .) a e a a i ab e f ee . f cha ge f . da a a a . .
 1 . e e a . . a ia da aba e ad a . ced faci 1 1 e [6].

I hi . . , e . ee . M . -SMOTI ha e e d SMOTI b . a i g ad a -
 age . f a igh 1 . e g a 1 . . i h a . a ia da aba e 1 . . de . . . i e . e 1 e
 a . . a ia . e ge . i . . . de f . . . i e a e . . The . . de 1 b 1 . a i g 1 . .
 acc . . a . h ee deg ee . f c . . e 1 . . . e e . ed ab . e .

3 Stepwise Mining of a Spatial Regression Model

M . -SMOTI . i e a . a ia . e ge . i . . . de b . a 1 1 . i g . a i g . a ia
 da a acc . di g . . 1 . a - a e . a d i e - a e . e a 1 . hi . a d a . cia i g dif
 fe . e . e ge . i . . . de . . . di 1 . . a ia a ea . I . a . ic a , i . . i e . a ia

da a a d e f . . . h e e i e c . . . c i . . . f a e e . . . c e d . . . d e i h
 b h i i g . . . d e a d e g e i . . . d e i . . . e . . . i g c i e i . . . i a i -
 e d . I h i a , i f a c e h e a i a e e d f d i g i h i g a . . . g e a a . .
 a i b e h a h a e . . . e g b a e e c . . . h e e . . . e a i b e a d h e h a
 h a e . . . c a e e c . B h i i g a d e g e i . . . d e . . . a i . . . e e e a
 a e a d a i a e a i . . . h i a . . . g h e .

Spatial split. A a i a . . . i i g e i . . . e e i h e a , . . . ,
 a , . e a e f . . . S . The f i e e a i i . . .
 a g e b e c . . . a c c d i g . . . e e a i a e a i i . . . h i (e i h e i a - a e . . . i e -
 a e). F i i a c e , h e e d i c i g h e i . . . f e e . . . e i g f . . .
 e i a . . . d i e a e i E D , i i a b e i g i c a . . . i e a d i e e . . . e g e i .
 f . c i . . . a c c d i g . . . h e e e c e . . . a b e c e f a i . . . a d c . . . i g h e e -
 A e a - c e . . . e c e f e f . . . i g . . . c h a i a e a i i . . . h i c . d i i .
 c . c e . . . h e i . . . d c i . . . f a . . . h e a e i . . . h e . . . d e . The a e i a e
 i . . . i g a b . . . e a c . d i i . ($X \leq \alpha$. . . $X > \alpha$) . . . h e c . . . i . . . c a e a d
 $X \in \{x_1, \dots, x_k\}$. . . $X \notin \{x_1, \dots, x_k\}$) . . . h e d i c e e . . . e . . . a h e a i c
 a i b e X f a a e a e a d i c d e d i h e . . . d e . I a d d i . . . h e a i c a -
 i b e , a a i b e c . d i i . . . a i . . . e a . . . a i a . . . e . . . (e.g., h e a e a f . . .
 . . . g . . . a d h e e e i . . . f . . . i e) , h a i i . . . i c i . . . d e d b h e g e . . . i -
 c a . . . c . . . e f h e c . . . e . . . d i g a e i S . I i i . . . e . . . h h a a i a
 e a i i . . . h i c . d i i . . . a d d e a e . . . f S . . . h e . . . d e . C h e e e . . . a i
 . . . a h e a i c a i b e . . . a i a . . . e . . . i . . . e e a a e a e a e a d i . . . d c e d
 i h e . . . d e . H e e e , d e . . . h e c . . . e i . . . f c . . . i g a i a e a i i . . . h i ,
 e i . . . e h a a e a i i . . . h i b e e e . . . a e . . . c a b e i . . . d c e d a c e
 i e a c h . . . i e a h c . . . e c i g h e h e e a f .

C h e e e . . . i h [10], h e a i d i . . . f a . . . a i a . . . i i g e i i b a e d . . . a
 h e i i c f . c i . . . $\sigma(t)$ h a i c e d . . . h e a i b e - a e e . . . e e a i . . .
 . . . f h e . . . i . . . f . . . a i a b e c . i S f a i g i t_L a d t_R , h a i , h e e f a d
 i g h c h i d f h e . . . i i g . . . d e t e e c i e . Thi a i b e - a e e . . . e e e -
 a i . . . c . . . e e . . . d . . . h e . . . e e f S d e i e d a c c d i g . . . b h . . . a i a e a -
 i . . . h i c . d i i . . . a d a i b e c . d i i . . . a . . . g h e a h f . . . h e f
 h e e e . . . h e c . . . e . . . d e . We d e e $\sigma(t) = (n(t_L)/(n(t_L) + n(t_R)))R(t_L) +$
 $(n(t_R)/(n(t_L) + n(t_R)))R(t_R)$, h e e $n(t_L)$ ($n(t_R)$) i h e . . . b e . . . f a i b e -
 a e . . . e e a e d d . . . h e e f (i g h) c h i d . S i c e i a - a e e a d i e - a e
 e a i i . . . h i e a d . . . a e g e i . . . d e h a . . . a i c d e e e a a e . . . (. . . e c -
 e a i . . . e a a e) , i i a h a e . . . h a $n(t) \neq n(t_L) + n(t_R)$ a h g h h e i
 i t a i e h e . . . a e c . i . . . e i e e . Thi i d e . . . h e a a
 . . . a e f i . . . a - a e e a d i e - a e e a i i . . . h i . I f a c , h e e e e a . . . a i a
 . . . b e c . a e . . . a i a e a e d . . . h e a e b e c (e.g., a i g e E D . . . a b e i e -
 . . . e c e d b e e e . . . a d) , c i g . . . a i a e a i i . . . h i . . . a e . . .
 a b e . . . f a i b e - a e . . . e g e a e h a . . . e . . . $R(t_L)$ ($R(t_R)$) i h e M i -
 . . . S a e d E . . . (MSE) c e d . . . h e e f (i g h) c h i d t_L (t_R) a f :

$$R(t_L) = \sqrt{\frac{1}{n(t_L)} \sum_{i=1..n(t_L)} (y_i - \hat{y}_i)^2} \quad (R(t_R) = \sqrt{\frac{1}{n(t_R)} \sum_{i=1..n(t_R)} (y_i - \hat{y}_i)^2},$$

ch ha \hat{g}_i he e... e a e , edic ed acc, di g... he... a ia , eg e...
 de b i b c... b i g he be... aigh -1 e, eg e... a... cia ed... $t_L(t_R)$,
 i h a... aigh -1 e, eg e... i... he a h f... he... $t_L(t_R)$ [3].

Spatial regression. A... a ia , eg e... de e f... a... aigh -1 e, eg e...
 ... er he a c... i... he a i c a... i b e... a c... i... a ia... e...
 ... e... d ced i... he... de c... e... b i... C he e... i h he e... e...
 ... ced... e [3], b... h e... e a d e... a a... a... i b e... a e... e... aced i h he...
 ... e i d a... F... i... a ce, he... a... eg e... i... e... i... e f... e d... a c... i...
 ... a... i b e X, he... e... e a... i b e... i... e... aced i h he... e i d a... $Y' = Y - Y$,
 ... he e... $Y = \alpha + \beta X$. The... eg e... i... c e... c i e... α a d β a e... e... i... a e d... he...
 a... i b e... a... e... e... e... a... i... f he... i... f S fa i g i... he c... e... de.

Acc, di g... he... a ia... c... e f... da a, he... eg e... i... a... i b e... c... e...
 f... e... e f he a e... a... e a d... i... e d... i... he... de... C... i... he... a i c a d...
 ... a ia... a... i b e... f he e... a e... , h i c h ha e... e... e... b e e... i... d ced i... he...
 ... de, a... e... e... aced i h he c... e... d i g... e i d a... i... d e... e... e... he...
 e... e c... f he... eg e... i... a... i b e... . W h e... e... a... e... a... e... i... a d d e... he... de...
 (b... e a... f a... a ia... e a... i... h i... c... d i...), c... i... he... a i c a d... a ia...
 a... i b e... , i... d ced i h i... , a... e... e... aced i h he c... e... d i g... e i d a... .
 R e i d a... a... e... c... e... a... c... e... d... he... a... i b e... a... e... e... e... a... i... f...
 he... i... f S fa i g i... he c... e... de. I... h i... a... , he... e... c... f... eg e... i...
 a... i b e... e... i... i... d ced i... he... de... b... eg e... i... e... i... a... e... e... d...
 b... i... d ced a... i b e... .

The e... a... a... i... f a... a ia... , eg e... i... e... $\hat{Y} = \hat{\alpha} + \hat{\beta}X$ i... b a e d... he...
 he... i... i c f... c... i... $\rho(t)$, ha... i... : $\rho(t) = \dots \{R(t), \sigma(t')\}$, he... e... t' i... he... be...
 ... a ia... i... i g... de f... i g... he... eg e... i... e... e... i... t. Th... e... -a h e a d... e...
 i... d e d... he... he... i... i c f... c... i... a b... e... d e... d... he... fac... ha... a ia... i...
 ... f... b e... aigh -1 e, eg e... i... a f e... he... i... c... d i... i... e... f... e d... h i e...
 he... eg e... i... e... d e... . A fa... e... c... a... i... d... b e... g... i g... he... e... a...
 a... f... he... e... e... b a... e... he... c... e... a... i... f... $\rho(T)$... he... be... i... e... i... e a...
 , eg e... i... a f e... he... eg e... i... e... e... X_i ... e... f... e d [10].

Stopping criteria. Th e... e d i... e... e... i g... c... i... e... i a... a... e... i... e... e... d... Th e...
 ... e... i... e... ha... a... i... a... b e... f... a... g e... b e c... f a... i... c... e... e... de. Th e... ec...
 ... d... e... he... i... d... c... i... c... e... he... he... c... e... c... i... e... f... d e... i... a... i... i... g... e a... e...
 ... ha... a... h e... h... d [18]. Th... c... e... c... i... e... i... a... c a... e... f... e... e... e... b e... e... a... f... he...
 ... e... g... h... f... he... e a... i... b e... e... e... e... a... a... a... i b e... i... he... ac... a... i... e...
 ... de... a... d... he... e... e... e... a... i b e... . F... i... a... , he... h... i... d... e... he... i... d... c... i... c... e...
 ... he... f... he... eg e... i... e... e... c a... b e... e... f... e d (i.e. a... c... i... e... a... i b e...
 a... e... i... c... d e d... i... he... c... e... e... de) a... a... f... e... i... d... c... i... g... e... e... e... a... e... .

4 Spatial Database Integration

M... a ia... da... a... i... i g... e... e... e... c e... da... a... i... a... i... e... e... . Th... e... e... i...
 h... i g... h... e... f... a... c... e... f... c... e... a... i... a... i... e... i... e... c... e... e... he... e... e... g... h... e...
 i... a... a... i b e... e... e... a... e... c... e... a... da... a... H... e... e... , i... a ia... da... a... i... e... i... e...
 ... c... e... e... i... i... i... a... e... e... i... e... i... e... f... e... e... c h a... i... f... f... a c... c... e... i g... ,... e... i g... a... d

1 de 1 g da a, cha h e a a r a b e 1 . . . a ia DBMS (Da a Ba e Ma age e . . . S e e . . .). F i a a ce, a ia e a 1 . . . (e.g., c . . . i g h e . . . g i c a e a . . . h i a . . . g a ia b e c . . .). . . . e d b a . . . a ia DBMS a e a d a age f . . . a ia 1 d e e i e Q a d , e e . . . Kd- . . . e e [14]. Thi . . . i a e a i g h 1 e- g a 1 . . . f a ia da a . . . i g . . . e . . . a d . . . a ia DBMS 1 . . . d e . . . 1) g a a e e h e a . . . i c a b i 1 . . . f a ia da a . . . i g a g , i h . . . a g e . . . a ia da a e . . . ; u) e . . . 1 . . . e f . . . e d g e f a ia da a . . . d e a r a b e , f e e f c h a g e , 1 h e . . . a ia da a b a e , m) . . . e c i f d i e c . . . h a da a . . . e d 1 a da a b a e h a e . . . b e . . . 1 e d , 1) a . . . i d . . . e e . . . e . . . c e i g e a d i g . . . e d . . . d a da a . . . a g e h a . . . a b e . . . e c e a . . . h e a . . . f a c e f h e h . . . h e i . . . a b e . . . e e . . . e . . . e d .

S . . . e e a . . . e . . . f i e g a i g . . . a ia da a . . . i g a d . . . a ia da a b a e . . . e . . . a e . . . e e . . . e d 1 [11] f . . . c a i c a 1 . . . a . . . a d 1 [1] f . . . a . . . c i a 1 . . . e d i c . . . e . . . a . . . I b h c a e , a da a . . . i g a g , i h i g 1 d e . . . g i c 1 e . . . 1 e g a e d i h a . . . a ia da a b a e b . . . e a . . . f . . . e . . . i d d e a e . . . d e h a e . . . a c . . . a ia a . . . i b e a d . . . e a 1 . . . h i . . . i d e e d e . . . f . . . h e . . . i g . . . e a d . . . e e . . . h e e f e a . . . e i a d e . . . g i c f . . . a - 1 Th . . . da a . . . i g a g , i h . . . a e . . . a c i c a . . . a . . . i e d . . . e . . . c e . . . e d da a a d h i . . . e . . . c e i g 1 . . . e - c e d . C . . . e . . . e , 1 [6] a . . . a ia da a . . . i g e . . . , a e d S b g . . . M i e , 1 e d f . . . h e a . . . f . . . b g . . . d i c . . . e . . . 1 . . . a ia da a b a e . S b g . . . d i c . . . e . . . 1 h e e a . . . a c h e d b a i g a d a - a g e f . . . a i g h 1 e g a 1 . . . f h e da a . . . i g a g , i h . . . i h h e da a b a e e . . . 1 . . . e . . . S a ia , e a 1 . . . h i . . . a d a . . . i b e a e h e d a i c a d e i e d b e . . . i 1 g . . . a ia DBMS e . . . e . . . 1 . . . f a c i 1 i e (e.g., a c a g e , c a . . . i d g e . . . e - e . . . d e . . .) a d . . . e d . . . g i d e h e . . . b g . . . d i c . . . e . . .

F . . . i g h e 1 . . . i a 1 . . . f S b g . . . M i e , e a . . . e a . . . b e c - e a 1 . . . a (OR) da a e . . . e e a 1 . . . , c h h a . . . a ia a e . . . e . . . e e i g b h . . . i 1 g a d e g e . . . 1 . . . d e a e e . . . e e d i h . . . a ia . . . e i e . The e . . . e i e i c d e . . . a ia . . . e a . . . b a e d . . . h e . . . - a . . . i c da a . . . e f . . . g e . . . e . . . c . . . i 1 g 1 a . . . d e . . . e d e . . . f c . . . d i a e (X,Y) . . . e . . . e e i g . . . 1 . . . , 1 e a d . . . g S i c e a ia . . . e a . . . 1 . . . e e . . . 1 b a i c e a 1 . . . a a g e b a . . . Da a g , e . . . e . . . a e e . . . 1 . . . f h e OR-DBMS O a c e S a ia C a . . . i d g e 9 h e e . . . a ia . . . e a c . . . e . . . a ia , e a 1 . . . h i . . . a d . . . e . . . a c . . . a ia a . . . i b e a e . . . a d e a r a b e f e e f c h a g e [6]. The e . . . e a . . . c a b e c a e d 1 SQL e . . . i e . . . F . . . e a . . . e : . . . *

Thi . . . a ia . . . e . . . e i e e h e a i . . . (ED, R a d) h e g i c a e a 1 . . . h i 1 d i . . . 1 = b . . . e a h e O a c e . . . e a R E L A T E . I 1 e . . . h h a , h e . . . e . . . f . . . c h SQL e . . . i e , a e a b e e d i e c a d . . . c h . . . e . . . a c i c a h a f a i g - . . . i a e e . . . 1 . . . f e a 1 . . . a a g e b a . . . Da a g . . . c h h a h e . . . e d i d i c a i da a b a e f a e [9].

W h e i g a . . . a ia . . . e . . . (a . . . c i a e d . . . a . . . d e f h e e e) , h e . . . e 1 a e . . . f . . . e d e c i b i g b h h e a i c a . . . i b e a d . . . a ia a . . . i b e . . . f 1 e d a e . . . The FROM c a . . . e i c d e a e . . . (. . . e c e a 1 d i e e . . .) 1 h e . . . d e a h e c e . . . d e . The WHERE c a . . . e i c d e . . . 1 c d i 1 . . . f . . . d a . . . g h e a h f . . . h e h e c e . . . d e . The e g a 1 . . . f e i h e a

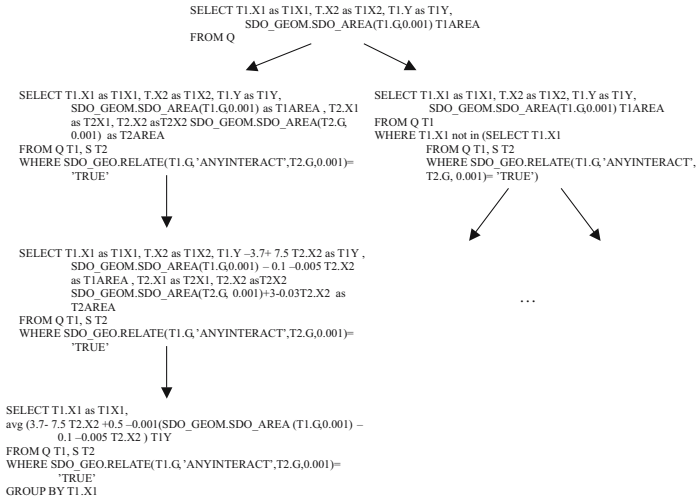


Fig. 2. An example of spatial model tree with regression, splitting and leaf nodes expressed by means of spatial queries assuming that training data are stored in spatial layers (e.g., Q and R) of a spatial database

... a ia , e a i . . . h i c . d i i . . . a a , i b e c . d i i . . . i g . . . e a , i b e
 f a . . . - a g e a e i , a f . . e d i . . a e g a e d . e e d . a i a . b - e .
 Thi i c h e e . . i h h e e a i c f e . . . i g . . i e a b e f a , e a -
 i . a d a b a e [2]. Fi a , he SELECT ca e i c d e h e a i c a d . a i a
 a , i b e (, h e i , e i d a .) f . . . h e a e . i . . . e d i h e WHERE ca e .
 Leaf . . d e a e a . . c i a e d i h a g g e g a i . . a i a . e i e , h a i , . . a i a
 e i e h e e a . . . e , e f e , i g h e a e a g e b e c a e g . . e d . g e h e .
 I h i a , h e , e d i c i . . f h e e . . . e a i a b e i h e a e a g e , e . . . e a e
 , e d i c e d . . h e e f a , i b e - a e . . e d e c i b i g h e . i e a g e b e c
 . b e , e d i c e d . Thi . e a . . h a . . a i a . . d e , e e c a b e e , e e d i f .
 f a . e . f S Q L . . a i a . e i e (e e F i g . 2) . Q e i e a e . . e d i X M L f . . a
 h a c a b e . b e e . . e d f . . e d i c i g (.) , e . . . e a , i b e .

5 Spatial Regression on Stockport Census Data

I h i e c i . . e . e e a , e a - , d a i c a i . . c e e i g h e . i i g f . a -
 i a , e g e . i . . . d e We c o u l d e b . h 1 9 9 1 c e . . . a d d i g i a . . a d a . . .
 . d e d i h e c . e . f h e E . . . e a . . . e c S P I N ! (S a i a M i . i g f . D a a f
 P b i c I . e e) [12]. Thi d a a c . c e . . S . c e f h e e . . e . . . i a
 d i . c i . G e a e , M a c h e e (U K) h i c h i d i d e d i . . e . - . . a d f .
 a . . a f 5 8 9 c e . . E D . S a i a a a . i i e a b e d b h e a a i a b i . . f e c -
 . . i e d b . . d a i e f , 5 7 8 S . c . . . E D a e e a b . h e , O . d a c e S . e
 d i g i a . . a . . f U K . D a a a e . . e d i a O . a c e S a i a C a . . i d g e Q d a b a e .

The attributes $h_1, d_1, e, g, a, e, h_2$ belong to the categorical domain. The attributes $ED, acc, d_1, g_1, h_1, b_1, g_1, a_1, e, e, a, i, a, b, e, f, e$ each belong to the numerical domain. The age group is a categorical variable with the values $ch, a, h, 1, g$ (53 years), $h, 1, g$ (9 years) and e, a, e, a (30 years) respectively. The ED variable has the values b, e, c, a, d respectively. The age group is a categorical variable with the values $ch, a, h, 1, g$ (53 years), $h, 1, g$ (9 years) and e, a, e, a (30 years) respectively.

The attributes $e, 1, e, a, e, 1, g, a, e, d, e, d$. The attributes $e, 1, g$ (BK₁) belong to the categorical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain.

The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain.

The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain.

The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain. The attributes $ED, The, ec, d, e, 1, g$ (BK₂) belong to the numerical domain.

¹ In both P1 and P2 transformations the attribute-value dataset is composed by 5 attributes for BK₁ (6 when including the lagged response) and 11 for BK₂ (12 when including the lagged response). The number of tuples for P1 is 4033 for BK₁ and 4297 for BK₂. In the case of P2, the number of tuples is 578 in both settings.

Table 1. Average MSE, No. of leaves and regression nodes of trees induced by Mrs-SMOTI, SMOTI and M5'. L1 is "No lagged response", L2 is "Lagged response".

Setting	MSE				Leaves				RegNodes				
	BK1		BK2		BK1		BK2		BK1		BK2		
	L1	L2	L1	L2	L1	L2	L1	L2	L1	L2	L1	L2	
Mrs-SMOTI	12.34	13.74	11.99	10.92	19.80	23.40	23.60	23.60	3.4	6.6	3.8	6.2	
SMOTI	P1	12.91	10.23	20.11	13.0	101.6	107.6	104.0	111.8	6.2	5.0	15.0	11.4
	P2	11.89	18.17	19.71	15.80	41.00	24.80	42.40	44.20	3.4	4.0	10.2	11.6
M5'	P1	13.52	12.41	12.92	12.30	433.6	872.0	408.6	711.2	-	-	-	-
	P2	12.44	9.19	12.48	9.59	198.0	199.4	199.2	197.4	-	-	-	-

Results show that Mrs-SMOTI1 has better accuracy than SMOTI1 in the classification task. This is due to the better handling of the regression nodes. The average MSE of the trees induced by SMOTI is significantly higher than the average MSE of the trees induced by Mrs-SMOTI (Table 1).

Moreover, the accuracy of SMOTI1 is significantly lower than the accuracy of Mrs-SMOTI1 (Table 1). This is due to the fact that Mrs-SMOTI1 has a higher accuracy than SMOTI1 in the classification task. The difference in accuracy between Mrs-SMOTI1 and SMOTI1 is statistically significant (Table 2). The average MSE of the trees induced by Mrs-SMOTI1 is significantly lower than the average MSE of the trees induced by SMOTI1 (Table 1). This is due to the fact that Mrs-SMOTI1 has a lower MSE than SMOTI1 in the regression task. The difference in MSE between Mrs-SMOTI1 and SMOTI1 is statistically significant (Table 2). The average MSE of the trees induced by Mrs-SMOTI1 is significantly lower than the average MSE of the trees induced by M5' (Table 1). This is due to the fact that Mrs-SMOTI1 has a lower MSE than M5' in the regression task. The difference in MSE between Mrs-SMOTI1 and M5' is statistically significant (Table 2).

Table 2. Mrs-SMOTI vs SMOTI and M5': results of the Wilcoxon test on the MSE of trees. If $W+ \leq W-$ then results are in favour of Mrs-SMOTI. The statistically significant values ($p \leq 0.1$) are in boldface. L1 is "No lagged response", L2 is "Lagged response".

Setting		Mrs-SMOTI vs. SMOTI P1			Mrs-SMOTI vs. SMOTI P2			Mrs-SMOTI vs. M5' P1			Mrs-SMOTI vs. M5' P2		
		W+	W-	p	W+	W-	p	W+	W-	p	W+	W-	p
		BK1	L1	6	9	0.81	9	6	0.81	3	12	0.310	7
	L2	10	5	0.63	6	9	0.81	8	7	1.000	15	0	0.060
BK2	L1	1	14	0.125	0	15	0.06	4	11	0.430	6	9	0.810
	L2	0	15	0.06	3	12	0.31	0	15	0.060	15	0	0.060

- **split** on *EDs'* number of migrants [≤ 47] (578 EDs)
 - **regression** on *EDs'* area (458 EDs)
 - **split** on *EDs - Shopping areas* spatial relationship (458 EDs)
 - **split** on *Shopping areas'* area (94 EDs) ...
 - **split** on *EDs'* number of migrants (364 EDs) ...
 - **split** on *EDs'* area (120 EDs)
 - **leaf** on *EDs'* area (22 EDs)
 - **regression** on *EDs'* area (98 EDs) ...

Fig. 3. Top-level description of a portion of the model mined by Mrs-SMOTI on the entire dataset at BK₂ level with no spatially lagged response attributes

The ... be ... f ... e ... 1 ... de ... a ... e ... a ... e ... 1 ... d ... a ... f ... he ... e ... 1 ... f ... he ... d ... ced ... e ... g ... e ... 1 ... de ... I ... h ... i ... ca ... e ... e ... h ... a ... h ... e ... de ... 1 ... d ... ced ... b ... M ... -SMOTI ... 1 ... ch ... i ... e ... h ... a ... h ... e ... de ... 1 ... d ... ced ... b ... SMOTI ... 1 ... b ... h ... e ... 1 ... g ... 1 ... de ... e ... de ... f ... da ... a ... a ... f ... a ... 1 ... The ... e ... a ... e ... 1 ... i ... c ... i ... f ... he ... a ... ia ... e ... g ... e ... 1 ... de ... 1 ... ed ... b ... M ... -SMOTI ... a ... e ... h ... e ... a ... 1 ... b ... e ... 1 ... e ... e ... ed ... I ... a ... i ... c ... a ... , ... h ... e ... e ... c ... e ... ca ... be ... e ... a ... 1 ... a ... g ... a ... ed ... 1 ... de ... 1 ... d ... i ... g ... 1 ... h ... a ... g ... g ... ba ... a ... d ... ca ... e ... ec ... f ... e ... a ... a ... a ... 1 ... b ... e ... F ... i ... a ... ce ... 1 ... Fig. 3 ... 1 ... h ... e ... - ... e ... de ... c ... 1 ... f ... he ... a ... ia ... e ... g ... e ... 1 ... de ... 1 ... ed ... b ... M ... -SMOTI ... h ... e ... 1 ... e ... da ... a ... e ... a ... BK₂ ... e ... 1 ... h ... a ... ia ... a ... g ... g ... e ... a ... 1 ... b ... e ... M ... -SMOTI ... ca ... e ... h ... e ... g ... ba ... e ... ec ... f ... he ... a ... e ... a ... f ... ED ... e ... S ... c ... e ... c ... e ... ed ... b ... he ... 458 ... ED ... h ... a ... 1 ... g ... e ... b ... e ... f ... i ... g ... a ... ≤ 47 ... The ... e ... ec ... f ... h ... e ... g ... e ... 1 ... 1 ... h ... a ... ed ... b ... a ... de ... 1 ... h ... e ... c ... e ... d ... i ... g ... b ... e ...

First, the ... a ... 1 ... f ... M ... -SMOTI ... 1 ... h ... M5', ... de ... h ... a ... ce ... a ... di ... e ... ce ... 1 ... e ... f ... MSE. ... a ... a ... M5', ... e ... e ... 1 ... a ... a ... di ... ad ... a ... g ... e ... 1 ... h ... e ... ec ... M ... -SMOTI. ... F ... i ... , ... M5' ... ca ... ca ... e ... a ... ia ... g ... ba ... a ... d ... ca ... e ... ec ... Sec ... d ... 1 ... ed ... de ... ce ... ca ... be ... 1 ... e ... e ... ed ... b ... h ... a ... be ... c ... e ... f ... he ... c ... e ... 1 ... f ... he ... de ... (... h ... e ... 1 ... a ... 1 ... c ... e ... a ... e ... f ... e ... de ... f ... a ... g ... 1 ... de ... 1 ... h ... e ... be ... f ... e ... a ... e ... f ... M ... -SMOTI ... M5')

6 Conclusions

I ... h ... i ... a ... e ... e ... h ... a ... e ... e ... e ... ed ... a ... a ... ia ... e ... g ... e ... 1 ... e ... h ... d ... M ... -SMOTI ... h ... a ... 1 ... a ... b ... e ... ca ... e ... b ... h ... a ... ia ... g ... ba ... a ... d ... ca ... e ... ec ... f ... e ... a ... a ... a ... 1 ... b ... e ... The ... e ... h ... d ... e ... ed ... h ... e ... 1 ... e ... c ... e ... c ... 1 ... f ... de ... ce ... e ... f ... ed ... b ... 1 ... e ... de ... ce ... SMOTI ... 1 ... d ... i ... ec ... 1 ... F ... i ... , ... b ... a ... 1 ... g ... ad ... a ... g ... e ... f ... a ... 1 ... h ... 1 ... e ... g ... a ... 1 ... 1 ... h ... a ... a ... ia ... da ... a ... b ... e ... de ... 1 ... 1 ... e ... b ... h ... a ... ia ... e ... a ... 1 ... h ... i ... a ... d ... a ... ia ... a ... 1 ... b ... e ... h ... i ... ch ... a ... e ... 1 ... i ... c ... i ... 1 ... a ... ia ... da ... a ... I ... de ... ed ... h ... i ... 1 ... i ... c ... i ... 1 ... f ... a ... 1 ... 1 ... f ... e ... e ... 1 ... b ... e ... f ... he ... a ... ia ... a ... ia ... 1 ... e ... de ... a ... a ... d ... i ... 1 ... e ... e ... e ... e ... ef ... 1 ... e ... g ... e ... 1 ... de ... 1 ... g ... Sec ... d ... , ... h ... e ... a ... ch ... a ... e ... g ... 1 ... d ... i ... ed ... 1 ... de ... 1 ... 1 ... e ... de ... h ... a ... ca ... e ... h ... e ... 1 ... i ... c ... e ... a ... 1 ... a ... c ... e ... f ... a ... ia ... da ... a ... Th ... e ... a ... h ... a ... ia ... e ... a ... 1 ... h ... i ... (1 ... a ... a ... e ... a ... d ... 1 ... e ... a ... e ...) ... a ... e ... 1 ... b ... e ... c ... o ... i ... de ... e ... a ... a ... a ... 1 ... b ... e ... h ... a ... 1 ... e ... ce ... h ... e ... e ... e ... a ... 1 ... b ... e ... b ... d ... e ... ce ... a ... 1 ... c ... e ... f ... a ... 1 ... g ... e ... a ... e ... I ... a ... i ... c ... a ... , ... 1 ... a ... a ... e ... e ... a ... 1 ... h ... i ... a ... e ... a ... ia ... g ... g ... e ... a ... 1 ... b ... e ... 1 ... a ... d ... i ... a ... ia ... g ... g ... e ... a ... 1 ... b ... e ... 1 ... a ... d ... i ... a ... ia ... g ... g ... e ... a ... 1 ... b ... e ...

Each element of the set of association rules is derived from the SMOTI. A few examples, related to the census data, are: $\{ \text{white} \} \rightarrow \{ \text{female} \}$ (e.g., 10%) and $\{ \text{white}, \text{female} \} \rightarrow \{ \text{married} \}$ (e.g., 15%).

Acknowledgment

This work is partially supported by the research project ATENEO-2005 of the University of Bari, Italy.

References

1. A. Appice, M. Ceci, A. Lanza, F. A. Lisi, and D. Malerba. Discovery of spatial association rules in georeferenced census data: A relational mining approach. *Intelligent Data Analysis*, 7(6):541–566, 2003.
2. A. Appice, M. Ceci, and D. Malerba. Mining model trees: A multi-relational approach. In T. Horvath and A. Yamamoto, editors, *Proceedings of ILP 2003*, volume 2835 of *LNAI*, pages 4–21. Springer-V., 2003.
3. N. R. Draper and H. Smith. *Applied regression analysis*. John Wiley & Sons, 1982.
4. S. Džeroski and N. Lavrač. *Relational Data Mining*. Springer-V., 2001.
5. R. Haining. *Spatial data analysis in the social and environmental sciences*. Cambridge University Press, 1990.
6. W. Klogsen and M. May. Spatial subgroup mining integrated in an object-relational spatial database. In T. Elomaa, H. Mannila, and H. Toivonen, editors, *Proceedings of PKDD 2002*, volume 2431 of *LNAI*, pages 275–286. Springer-V., 2002.
7. J. Knobbe, M. Haas, and A. Siebes. Propositionalisation and aggregates. In L. D. Raedt and A. Siebes, editors, *Proceedings of PKDD 2001*, volume 2168 of *LNAI*, pages 277–288. Springer-V., 2001.
8. K. Koperski. *Progressive Refinement Approach to Spatial Data Mining*. PhD thesis, Computing Science, Simon Fraser University, British Columbia, Canada, 1999.
9. G. Kuper, L. Libkin, and L. Paredaens. *Constraint databases*. Springer-V., 2001.
10. D. Malerba, F. Esposito, M. Ceci, and A. Appice. Top down induction of model trees with regression and splitting nodes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(5):612–625, 2004.
11. D. Malerba, F. Esposito, A. Lanza, F. A. Lisi, and A. Appice. Empowering a gis with inductive learning capabilities: The case of ingens. *Journal of Computers, Environment and Urban Systems, Elsevier Science*, 27:265–281, 2003.
12. M. Wang. Spatial knowledge discovery: The spin! system. In K. Fullerton, editor, *Proceedings of the EC-GIS Workshop*, 2000.
13. M. Orkin and R. Drogin. *Vital Statistics*. McGraw Hill, New York, USA, 1990.
14. H. Samet. *Applications of spatial data structures*. Addison-Wesley longman, 1990.
15. S. Shekhar, P. R. Schrater, R. Vatsavai, W. Wu, and S. Chawla. Spatial contextual classification and prediction models for mining geospatial data. *IEEE Transactions on Multimedia*, 4(2):174–188, 2002.
16. L. Torgo. *Inductive Learning of Tree-based Regression Models*. PhD thesis, Department of Computer Science, University of Porto, Porto, Portugal, 1999.
17. Y. Wang and I. Witten. Inducing model trees for continuous classes. In M. Van Someren and G. Widmer, editors, *Proceedings of ECML 1997*, pages 128–137, 1997.
18. S. Weisberg. *Applied regression analysis*. Wiley, New York, USA, 1985.

Word Sense Disambiguation for Exploiting Hierarchical Thesauri in Text Classification

Dimitrios Mavroeidis¹, George Tsatsaronis¹, Michalis Vazirgiannis¹,
Martin Theobald², and Gerhard Weikum²

¹Department of Informatics, Athens University of Economics and Business, Greece

²Max-Planck Institute of Computer Science,
Saarbruecken, Germany

Abstract. The introduction of hierarchical thesauri (HT) that contain significant semantic information, has led researchers to investigate their potential for improving performance of the text classification task, extending the traditional “bag of words” representation, incorporating syntactic and semantic relationships among words. In this paper we address this problem by proposing a Word Sense Disambiguation (WSD) approach based on the intuition that word proximity in the document implies proximity also in the HT graph. We argue that the high precision exhibited by our WSD algorithm in various humanly-disambiguated benchmark datasets, is appropriate for the classification task. Moreover, we define a semantic kernel, based on the general concept of GVSM kernels, that captures the semantic relations contained in the hierarchical thesaurus. Finally, we conduct experiments using various corpora achieving a systematic improvement in classification accuracy using the SVM algorithm, especially when the training set is small.

1 Introduction

It can be argued that WSD algorithms for the document classification task should differ in their design and evaluation from pure WSD algorithms. It is expected that correctly disambiguated words could improve (and certainly not degrade) the performance of a document classification task, while falsely disambiguated words would entail noise. Although the SVM algorithm [1] used in our experiments is known to be noise tolerant, it is certain that noise, above a certain level, will eventually degrade in SVM’s performance. In the absence of theoretical or experimental studies on the exact level of falsely disambiguated words that can be tolerated by classification algorithms, the most appropriate performance measure for WSD algorithms designed for a classification task is precision. Choosing the WSD algorithm with the highest precision will result in the incorporation of the lowest amount of noise in the classification task.

Another important issue for the successful embedding of WSD in text classification, is the exploitation of senses’ semantic relations, that are provided by the HT. These relations are essential for defining distances and kernels that reflect semantic similarities between senses. An extensive bibliography exists for measuring distances and similarities on thesauri and ontologies, which has not been taken into account by other research

approaches that embed WSD in the text classification task. The need for exploiting semantic relations is illustrated in [2], where SemCor 1.7.1, a humanly-disambiguated corpus, is used in classification experiments. It is demonstrated that even with a 100% accurate disambiguation, the simple use of senses instead of keywords does not improve classification performance.

In this paper we propose an unsupervised WSD algorithm for classification, that utilizes a background HT. Our approach adopts the intuition that adjacent terms extracted from a given document are expected to be semantically close to each other and that is reflected to their pathwise distance on the HT. Thus, the objective of our WSD method is, given a set of terms, to select the senses (one for each term among many found in the HT) that overall minimize the pathwise distance and reflect the compactness of the selected sense set. The semantic compactness measure introduced is based on the concept of the Steiner Tree [3]. As opposed to other approaches that have utilized WSD for classification [4],[5],[6],[7], we have conducted extensive experiments with disambiguated corpora (Senseval 2 and 3, SemCor 1.7.1), in order to validate the appropriateness of our WSD algorithm. Experiments, using the WordNet HT, demonstrate that our WSD algorithm can be configured to exhibit very high precision, and thus can be considered appropriate for classification. In order to exploit the semantic relations inherent in the HT, we define a semantic kernel based on the general concept of GVSM kernels [8]. Finally, we have conducted experiments utilizing various sizes of training sets for the two largest Reuters-21578 categories and a corpus constructed from crawling editorial reviews of books from the Amazon website. The results demonstrate that our approach for exploiting hierarchical thesauri semantic information contributes significantly to the SVM classifier performance, especially when the training set size is small.

In the context of this paper WordNet [9] is utilized as a hierarchical thesaurus both for WSD and for classification. Although WordNet contains various semantic relations between concepts¹, our approach relies only on the hypernym/hyponym relation that orders concepts according to generality, and thus our approach can generalize to any HT that supports the hypernym/hyponym relation.

The rest of the paper is organized as follows. Section 2 discusses the preliminary notions and the related work. Section 3 presents our compactness measure for WSD that is based on the graph structure of an HT. Section 4 describes the semantic kernel that is utilized for the experiments. Section 5 discusses the experiments performed. Section 6 contains the comparison of the proposed framework to other approaches, concluding remarks and pointers to further work.

2 Preliminaries

2.1 Graph Theoretic Notions

Assuming that a document is represented by a set of senses, the semantic compactness measure that we introduce for WSD implies a similarity notion either among the senses of a sense set or between two sense sets. Its commutation is based on the notion of

¹ Concepts are word senses in WordNet terminology and in this paper we will use the terms word senses and concepts interchangeably.

Steiner Tree. Given a set of graph vertices, the Steiner Tree is the smallest tree that connects the set of nodes in the graph. The formal definition of the Steiner Tree is given below.

Definition 1 (Steiner Tree). *Given an undirected graph $G = (V, E)$, and a set $S \subseteq V$, then the Steiner Tree is the minimal Tree of G that contains all vertices of S .*

2.2 Semantic Kernels Based on Hierarchical Thesaurus

Since we aim at embedding WSD in the SVM classifier, we require the definition of a kernel that captures the semantic relations provided by the HT. To the extend of our knowledge the only approach that defines a semantic kernel based on a HT is [10]. The formal definition of their kernel is given below.

Definition 2 (Semantic Smoothing Kernels [10]). *The Semantic smoothing Kernel between two documents d_1, d_2 is defined as $K(d_1, d_2) = d_1 P' P d_2 = d_1 P^2 d_2$, where P is a matrix whose entries $P_{ij} = P_{ji}$, represent the semantic proximity between concepts i and j .*

The similarity matrix P is considered to be derived by a HT similarity measure. The Semantic Smoothing Kernels have similar semantics to the GVSM model defined in [8]. A kernel definition based on the GVSM model is given below.

Definition 3 (GVSM Kernel). *The GVSM kernel between two documents d_1 and d_2 is defined as $K(d_1, d_2) = d_1 D D' d_2$, where D is the term document matrix.*

The rows of matrix D , in the GVSM kernel contain the vector representation of terms, used to measure their pairwise semantic relatedness. The Semantic Smoothing Kernel has similar semantics. The Semantic Smoothing Kernel between two documents $K(d_1, d_2) = d_1 P^2 d_2$, can be regarded as a GVSM kernel, where the matrix D is derived by the decomposition of $P^2 = D D'$ (the decomposition is always possible since P^2 is guaranteed to be positive definite). The rows of D can be considered as the vector representation of concepts, used to measure their semantic proximity. Semantic Smoothing Kernels use P^2 and not P , because P is not guaranteed to be positive definite.

2.3 Related Work

WSD. The WordNet HT has been used for many supervised and unsupervised WSD algorithms. In direct comparison to our WSD approach we can find [11],[12],[13] that are unsupervised and rely on the semantic relations provided by WordNet. In the experimental section we show that our WSD algorithm can be configured to exhibit very high precision in various humanly-disambiguated benchmark corpora, and thus is more appropriate for the classification task.

Senseval (www.senseval.org), provides a forum, where the state of the art WSD systems are evaluated against disambiguated datasets. In the experimental sections we will compare our approach to the state of the art systems that have been submitted to the Senseval contests.

WSD and classification. In this section we shall briefly describe the relevant work done in embedding WSD in the document classification task. In [7], a WSD algorithm based on the general concept of Extended Gloss Overlaps is used and classification is performed with an SVM classifier for the two largest categories of the Reuters-25178 collection and two IMDB movie genres (www.imdb.com).

It is demonstrated that, when the training set is small, the use of WordNet senses together with words improves the performance of the SVM classification algorithm, however for training sets above a certain size, the approach is shown to have inferior performance to term-based classification. Moreover, the semantic relations inherent in WordNet are not exploited in the classification process. Although the WSD algorithm that is employed is not verified experimentally, its precision is estimated with a reference to [13], since the later work has a very similar theoretical basis. The experiments conducted by [13] in Senseval 2 lexical sample data, show that the algorithm exhibits low precision (around 45%) and thus may result in the introduction of much noise that can jeopardize the performance of a classification task.

In [4], the authors experiment with various settings for mapping words to senses (no disambiguation, most frequent sense as provided by WordNet and WSD based on context). Their approach is evaluated on the Reuters-25178, the OSHUMED and the FAODOC corpus, providing positive results. Their WSD algorithm has similar semantics to the WSD algorithm proposed in [12]. Although in [12] the experiments are conducted in a very restricted subset of SemCor 1.7.1, the results reported can be compared with our experiment results for the same task, as it is shown in Section 5. Moreover [4], use hypernyms for expanding the feature space.

In [5] the authors utilize the supervised WSD algorithm proposed in [14] in k-NN classification of the 20-newsgroups dataset. The WSD algorithm they employ is based on a Hidden Markov Model and is evaluated against Senseval 2, using “English all words task”, reporting a maximum precision of around 60%. On the classification task of the 20-newsgroup dataset, they report a very slight improvement in the error-percentage of the classification algorithm. The semantic relations that are contained in WordNet are not exploited in the k-NN classification process.

The authors in [6] present an early attempt to incorporate semantics by means of a hierarchical thesauri in the classification process, reporting negative results on the Reuters-21578 and DigiTrad collection. While none disambiguation algorithm is employed, the use of hypernyms for extending the feature space representation is levied.

2.4 Hierarchical Thesaurus Distances – Similarities

As we have discussed in the introduction section, an important element for the successful incorporation of semantics in the classification process is the exploitation of the vast amount of semantic relations that are contained in the HT. There is an extensive bibliography that addresses the issue of defining distances and similarity measures based on the semantic relations provided by an HT [9],[15],[16],[17], which has not been related to the existing approaches for embedding WSD in classification. A common ground of most of the approaches is that the distance or similarity measure will depend on the “size” of the shortest path that connects the two concepts through a common ancestor in the hierarchy, or on the largest “depth” of a common ancestor in the hierarchy. The

terms “size” and “depth” are used in an informal manner, for details one should use the references provided.

3 Compactness Based Disambiguation

In this section we present our unsupervised WSD method, as this was initially sketched in [18]. Our WSD algorithm is based on the intuition that adjacent terms extracted from a text document are expected to be semantically close to each other. Given a set of adjacent terms, our disambiguation algorithm will consider all the candidate sets of senses and output the set of senses that exhibits the highest level of semantic relatedness. Therefore, the main component of our WSD algorithm is the definition of a semantic compactness measure for sets of senses. We refer to our disambiguation approach as CoBD (Compactness Based Disambiguation). The compactness measure utilized in CoBD is defined below.

Definition 4. *Given an HT O and a set of senses $S = (s_1, \dots, s_n)$, where $s_i \in O$ the compactness of S is defined as the cost of the Steiner Tree of $S \cup lca(S)$, such that there exists at least one path, using hypernym relation, from each s_i to the $lca(S)$.*

In the definition above we include one path, using the hypernym relation, for every sense to the least common ancestor $lca(S)$. The reason for imposing such a restriction is that the distance between two concepts in an HT is not defined as the shortest path that connects them in the HT, but rather as the shortest path that goes through a common ancestor. Thus, it can be argued that two concepts are connected only through a common ancestor and not through any other path in the HT. The existence of the $lca(S)$ (and of a path between every concept and the $lca(S)$ using the hypernym relation) guarantees that a path connecting all pairs of concepts (in the context discussed earlier) exists.

Although in general the problem of computing the Steiner Tree is NP-complete, the computation of the Steiner Tree (with the restriction imposed) of a set of concepts with their lca in a HT is computationally feasible and is reduced to the computation of the shortest path of the lca to every concept of the set. Another issue, potentially adding excessive computational load, is the large number of combinations of possible sets of senses, when a term set of large cardinality is considered for disambiguation. In order to address this issue, we reduce the search space by using a Simulated Annealing algorithm. The experimental setup used in this paper for the empirical evaluation of our WSD algorithm is described in detail in section 5.

4 Exploitation of Hierarchical Thesaurus Semantics in SVM Classification

We have argued in the introductory section that the exploitation of the semantics provided by an HT are important for the successful embedding of WSD in the classification task. In this section we will present the definition of the Kernel we will utilize in SVM classification. The Kernel we define is based on the general concept of GVSM kernel and depicts the semantics of the HT.

It is shown in detail in [18], that the use of hypernyms for the vector space representation of the concepts of a HT, enables the measurement of semantic distances in the vector space. More precisely, given a Tree HT, there exists a weight configuration for the hypernyms, such that standard vector space distance and similarity measures are equivalent to popular HT distances and similarities. The proofs for propositions given below can be found in [18].

Proposition 1. *Let O be a Tree HT, if we represent the concepts of the HT O , as vectors containing all their hypernyms, then there exists a configuration for the weights of the hypernyms such that the Manhattan distance (Minkowski distance with $p=1$) of any two concepts in vector space is equal to the Jiang-Conrath measure [15] in the HT.*

Proposition 2. *Let O be a Tree HT, if we represent the concepts of the HT O as vectors containing all their hypernyms, then there exists a configuration for the weights of the hypernyms such that the inner product of any two concepts in vector space is equal to the Resnik similarity measure [16] in the HT.*

The WordNet hierarchical thesaurus is composed by 9 hierarchies that contain concepts that inherit from more than one concept, and thus are not Trees. However, since only 2.28% of the concepts inherit from more than one concept [19], we can consider that the structure of WordNet hierarchies is close to the Tree structure.

From the above we conclude that, if we construct a matrix D where each row contains the vector representation of each sense containing all its hypernyms, the matrix DD' will reflect the semantic similarities that are contained in the HT. Based on D , we move on to define the kernel between two documents d_1, d_2 , based on the general concept of GVSM kernels as $K_{concepts}(d_1, d_2) = d_1 DD' d_2$. In our experiments we have used various configurations for the rows of D . More precisely, we have considered the vector representation of each concept to be extended with a varying number of hypernyms. The argument for using only a limited number and not all hypernyms is that the similarity between hypernyms close to the root of the HT is considered to be very close to 0. Apart from hypernyms, in the experiments section, we have explored the potential of using hyponyms for constructing matrix the D in the GVSM kernel. The kernel that we finally utilize in our experiments is a combination of the inner product kernel for terms with the concept kernel $K(d_1, d_2) = K_{terms}(d_1, d_2) + K_{concepts}(d_1, d_2)$. This kernel was embedded into the current version of *SVMLight* [20] and replaced the standard linear kernel used for document classification with sparse training vectors.

The kernel defined implies a mapping from the original term and concept space, to a space that includes the terms, the concepts and their hypernyms. The kernel can be considered as the inner product in this feature space.

5 Experiments

5.1 Evaluation of the WSD Method

CoDB was tested in four benchmark WSD corpora; Brown 1 and Brown 2 from the SemCor 1.7.1 corpus, and the in the “English All Words” task of Senseval 2 and 3. These corpora are pre-tagged and pre-annotated. From all the parts of speech in the

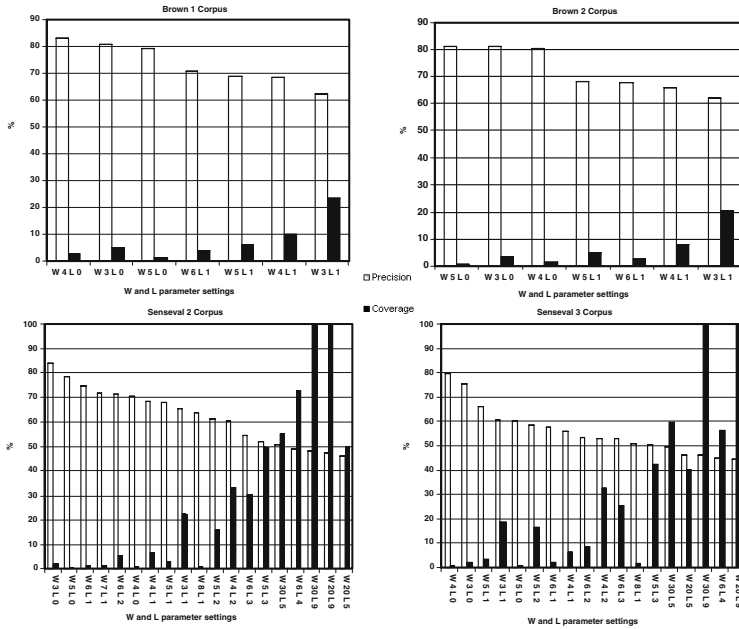


Fig. 1. WSD results on 4 benchmark datasets for different initializations of W and L

texts we only considered nouns, which are usually more informative than the rest and form a meaningful type hierarchy in WordNet. In order to implement CoBD efficiently we had to take into account that the search space of combinations to be examined for their compactness increases dramatically as the cardinality of the set of words examined increases, making exhaustive computation infeasible. Thus we adopted simulated annealing as in [21]. This approach reduced the search space and allowed us to execute the WSD using various set of words sizes in a time efficient manner. The parameters of the WSD method are:

1. Window Size (W): Set cardinality of the words to be disambiguated.
2. Allowed Lonely: Given a word set L , it is the maximum number of lonely senses¹ allowed in a WordNet noun hierarchy, for any senses combination of that window.

Figure 1 presents experiments we have conducted using various parameter settings. The results are sorted in decreasing order of precision. The precision and coverage² values reported do not take into account the monosemous nouns, but only the ambiguous ones. We can estimate, based on the examined corpora statistics, that the inclusion of the monosemous nouns would report an increase in precision between 3% and 4%, as well as an increase in coverage of almost 22%.

¹ A sense s belonging to a set of senses S is referred to as lonely if the WordNet noun hierarchy H it belongs to, does not contain any other $k \in S$.

² Coverage is defined as the percentage of the nouns that are disambiguated.

We observe that CoBD achieves precision greater than 80% with an associated coverage of more than 25%, if monosemous (i.e., non-ambiguous) nouns are also taken into account. Comparable experiments conducted in [12] reported a top precision result of 64,5% with an associated coverage of 86,2%. Similar experiments conducted in [11], [13] and [14] resulted as well in lower precision than CoBD. In comparing our approach to the state of the art WSD algorithms that were submitted to the “English All Words” Senseval 2 contest (www.senseval.org), we observe that our approach can be configured to exhibit the highest precision.

5.2 Document Collections and Preprocessing for Text Classification

Reuters. Reuters-21578 is a compilation of news articles from the Reuters newswire in 1987. We include this collection mostly for transparency reasons, since it has become the gold standard in document classification experiments. We conducted experiments on the two largest categories, namely *acquisitions* and *earnings*, in terms of using test- and training documents based on the [4] split. This split yields a total of 4,436 training and 1,779 test documents for the two categories. We extracted features from the mere article bodies, thus using whole sentences only and hiding any direct hint to the actual topic from the classifier.

Amazon. To test our methods on a collection with a richer vocabulary, we also extracted a real-life collection of natural-language text from amazon.com using Amazon’s publicly available Web Service interface. From that taxonomy, we selected all the available editorial reviews for books in the three categories *Physics*, *Mathematics* and *Biological Sciences*, with a total of 6,167 documents. These reviews typically contain a brief discussion of a book’s content and its rating. Since there is a high overlap among these topics’ vocabulary and a higher diversity of terms within each topic than in Reuters, we expect this task to be more challenging for both the text- as well as the concept-aware classifier.

Before actually parsing the documents, we POS-annotated both the Reuters and Amazon collections, using a version of the commercial Connexor software for NLP processing. We restricted the disambiguation step to matching noun phrases in WordNet, because only noun phrases form a sufficiently meaningful HT in the ontology DAG. Since WordNet also contains the POS information for each of its concepts, POS document tagging significantly reduces the amount of choices for ambiguous terms and simplifies the disambiguation step. For example the term *run* has 52 (!) distinct senses in WordNet out of which 41 are tagged as verbs. The parser first conducts continuous noun phrase tokens in a small window of up to a size of 5 into dictionary lookups in WordNet before the disambiguation step takes place. If no matching phrase is found within the current window, the window is moved one token ahead. This sliding window technique enables us to match any composite noun phrase known in WordNet, whereupon larger phrases are typically less ambiguous. Non-ambiguous terms can be chosen directly as safe seeds for the compactness-based disambiguation step. Note that we did not perform any feature selection methods such as Mutual Information or Information Gain [22] prior to training the SVM, in order not to bias results toward a specific classification method.

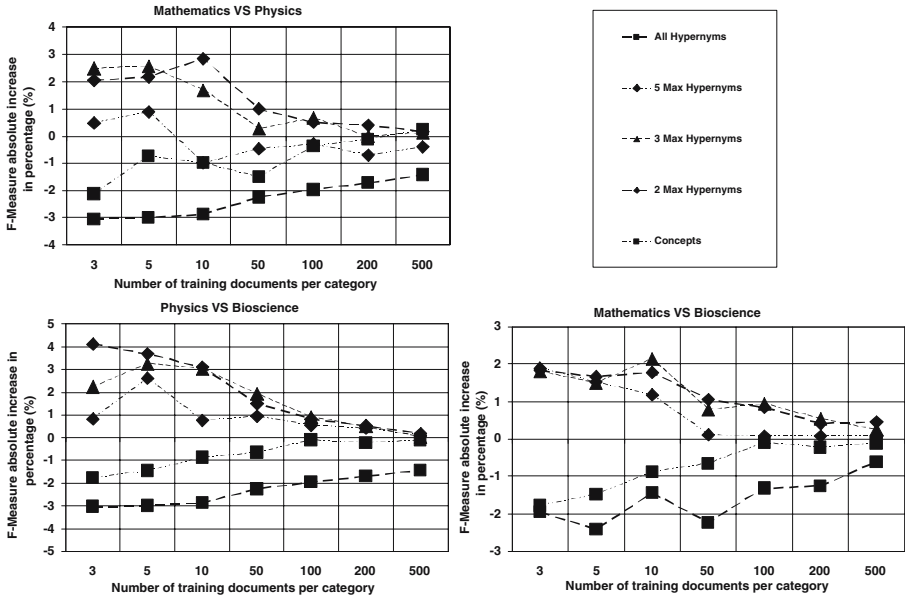


Fig. 2. Relative Improvement of F-measures scores for various Similarity Configurations in the Amazon Topics

5.3 Evaluation of Embedding CoBD in the Text Classification Task

To evaluate the embedding of CoBD in text classification, we performed binary classification tasks, only, i.e., we did not introduce any additional bias from mapping multi-class classification task onto the binary decision model used by the SVM method. The binary classification tasks were performed after forming all pairs between the three Amazon topics, and one pair between the two largest Reuters-21578 topics. The parameters' setting for CoBD was $W3L0$, since it reported high precision and performed in a stable manner during the WSD evaluation experiments in the 4 benchmark corpora. Our baseline was the F-Measure [22] arising from the mere usage of term features. The baseline competed against the embedding of the term senses, whenever disambiguation was possible, and their hypernyms/hyponyms into the term feature vectors, according to the different GVSM kernel configurations shown in Figures 2,3. In our experiments, the weights of the hypernyms used in the GVSM kernel are taken to be equal to the weights of the terms they correspond to. We varied the training set sizes between 3 and 500 documents per topic. For each setup, in Figures 2,3 we report the differences of the *macro-averaged F-Measure* between the baseline and the respective configurations, using 10 iterations for each of the training set sizes of the Reuters dataset and 30 iterations for each of the training set sizes of the Amazon dataset. The variation of the differences was not too high and allowed for all the results where the absolute difference of the sample means was greater than 1% to reject the null hypothesis (that the means are equal) at a significance level of 0.05. For more than 500 documents, all our experiments

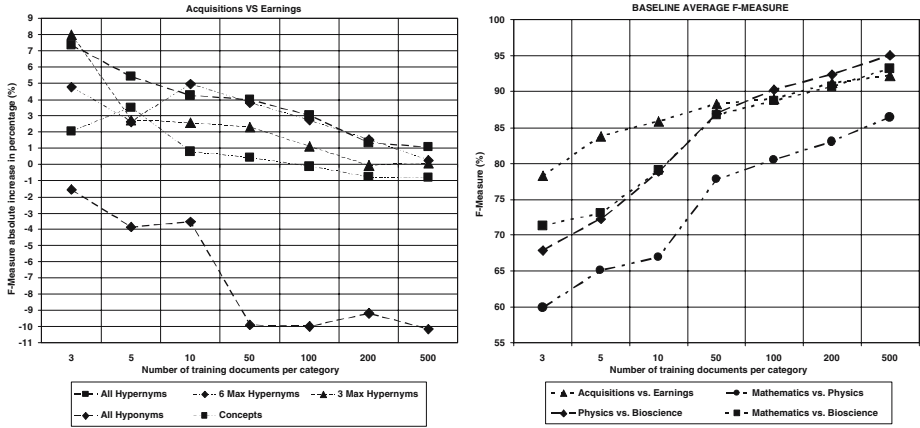


Fig. 3. Relative Improvement of F-measures scores for various Similarity Configurations in the Reuters Topics

indicate a convergence in results between the concept-aware classifier and the text classifier. The average F-measures for the baseline classifier are reported in Figure 3. For each run, the training documents were selected randomly following a uniform distribution. Since there is no split into separate documents for training and testing given in the Amazon collection, we performed cross-validation runs over the whole set, each using all the remaining documents for the test phase.

The results demonstrate that the use of CoBD and our kernel function, based on a small number of hypernyms increases consistently the classification quality especially for small training sets. In some cases, as the number of hypernyms increases we observe a performance deterioration which in some cases falls below the term-based classification.

The variance in the number of hypernyms needed for achieving better performance, can be explained by the fact that we did not employ a hypernym weighting scheme. Thus, when semantically correlated categories are considered, (such as Maths/Physics in the Amazon data), then the use of all the hypernyms with equal weights would result in many documents belonging to the Physics category to have a high similarity to documents of Maths category, degrading the performance of the classification algorithm.

6 Discussion and Conclusions

The context of the current work entails the content and structure (i.e. the senses and hierarchical relationships) of HTs and their usage for successful extension of the bag of words model for text classification. The objective is that such extensions (i.e. senses and hypernyms/hyponyms more precisely) are contributing to higher quality in the classification process.

The *contribution* of the paper is the design of a successful WSD approach to be incorporated and improve the text classification process. Our WSD approach takes into account term senses found in HTs, (in the specific case Wordnet), and for each document

selects the best combination of them based on their conceptual compactness in terms of related Steiner tree costs. Apart from the senses we add to the original document feature set a controlled number of hypernyms of the senses at hand. The hypernyms are incorporated by means of the kernel utilized. The attractive features of our work are:

Appropriate WSD approach for text classification. Most of the related approaches incorporating WSD in the classification task [6],[7],[4] do not provide a sound experimental evidence on the quality of their WSD approach. On the contrary in our work, the WSD algorithm is exhaustively evaluated against various humanly disambiguated benchmark datasets and achieves very high precision (among the top found in related work) although at low coverage values (see Fig.1). This is not a problem, though since as mentioned earlier, it is essential to extend the feature space with correct features in order to prevent introduction of noise in the classification process. The experimental evaluation provides us with the assurance that our WSD algorithm can be configured to have high precision, and thus, would insert in the training set very little noise.

Similarity measure that takes into account the structure of the HT. Document classification depends on a relevant similarity measure to classify a document into the closest of the available classes. It is obvious that the similarity among sets of features (representing documents) should take into account their hierarchical relationships as they are represented in the HT. None of the previous approaches for embedding WSD in classification has taken into account the existing literature for exploiting the HT relations. Even when the use of hypernyms is used [6],[4], it is done in an ad-hoc way, based on the argument that the expansion of a concept with hypernyms would behave similar to query expansion using more general concepts. We utilize a Kernel based on the general concept of a GVSM kernel that can be used for measuring the semantic similarity between two documents. The kernel is based on the use of hypernyms for the representation of concepts - theoretically justified in the context of the related work concerning the computation of semantic distances and similarities on a HT that aligns to tree structure.

We conducted classification experiments on two real world datasets (the two largest Reuters categories and a dataset constructed by the editorial reviews of products on three categories at the *amazon.com* web site). The results demonstrate that our approach for embedding WSD in classification yields significantly better results especially when the training sets are small.

An issue that we will investigate in further work is the introduction of a weighting scheme for hypernyms favoring hypernyms that are close to the concept. A successful weighting scheme is expected to reduce the problem of the variance in the number of hypernyms needed to achieve optimal performance. We will investigate learning approaches to learn the weighting schemes for hypernyms. Moreover, we aim in conducting further experiments on other larger scale and heterogeneous data sets.

References

1. Drucker, H., Burges, C.J.C., Kaufman, L., Smola, A.J., Vapnik, V.: Support vector regression machines. In: Advances in Neural Information Processing Systems (NIPS). (1996) 155–161
2. Kehagias, A., Petridis, V., Kaburlasos, V.G., Fragkou, P.: A comparison of word- and sense-based text categorization using several classification algorithms. *Journal of Intelligent Information Systems* **21** (2003) 227–247

3. Hwang, R., Richards, D., Winter, P.: The steiner tree problem. *Annals of Discrete Mathematics* **53** (1992)
4. Bloehdorn, S., Hotho, A.: Boosting for text classification with semantic features. In: Proc. of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Mining for and from the Semantic Web Workshop. (2004) 70–87
5. Rosso, P., Ferretti, E., Jimenez, D., Vidal, V.: Text categorization and information retrieval using wordnet senses. In: Proc. of the 2nd International WordNet Conference (GWC). (2004)
6. Scott, S., Matwin, S.: Feature engineering for text classification. In: Proc. of the 16th International Conference on Machine Learning (ICML). (1999) 379–388
7. Theobald, M., Schenkel, R., Weikum, G.: Exploiting structure, annotation, and ontological knowledge for automatic classification of xml data. In: International Workshop on Web and Databases (WebDB). (2003) 1–6
8. Wong, S.K.M., Ziarko, W., Wong, P.C.N.: Generalized vector space model in information retrieval. In: Proc. of the 8th annual international ACM SIGIR conference on Research and development in information retrieval. (1985) 18–25
9. Fellbaum, C., ed.: *WordNet, An Electronic Lexical Database*. The MIT Press (1998)
10. Siolas, G., d’Alche Buc, F.: Support vector machines based on semantic kernel for text categorization. In: Proc. of the International Joint Conference on Neural Networks (IJCNN). Volume 5., IEEE Press (2000) 205–209
11. Sussna, M.: Word sense disambiguation for free-text indexing using a massive semantic network. In: Proc. of the 2nd International Conference on Information and Knowledge Management (CIKM). (1993) 67–74
12. Agirre, E., Rigau, G.: A proposal for word sense disambiguation using conceptual distance. In: Proc. of Recent Advances in NLP (RANLP). (1995) 258–264
13. Banerjee, S., Pedersen, T.: Extended gloss overlaps as a measure of semantic relatedness. In: Proc. of the 18th International Joint Conference on Artificial Intelligence (IJCAI). (2003) 805–810
14. Molina, A., Pla, F., Segarra, E.: A hidden markov model approach to word sense disambiguation. In: Proc. of the 8th Iberoamerican Conference on Artificial Intelligence. (2002)
15. Jiang, J., Conrath, D.: Semantic similarity based on corpus statistics and lexical taxonomy. In: Proc. of the International Conference on Research in Computational Linguistics. (1997)
16. Resnik, P.: Using information content to evaluate semantic similarity in a taxonomy. In: Proc. of the 14th International Joint Conference on Artificial Intelligence (IJCAI). (1995)
17. Lin, D.: An information-theoretic definition of similarity. In: Proc. of the 15th International Conference on Machine Learning (ICML). (1998) 296–304
18. Mavroeidis, D., Tsatsaronis, G., Vazirgiannis, M.: Semantic distances for sets of senses and applications in word sense disambiguation. In: Proc. of the 3rd International Workshop on Text Mining and its Applications. (2004)
19. Devitt, A., Vogel, C.: The topology of wordnet: Some metrics. In: Proc. of the 2nd International WordNet Conference (GWC). (2004) 106–111
20. Klinkenberg, R., Joachims, T.: Detecting concept drift with support vector machines. In: Proc. of the 17th International Conference on Machine Learning (ICML). (2000) 487–494
21. Cowie, J., Guthrie, J., Guthrie, L.: Lexical disambiguation using simulated annealing. In: 14th International Conference on Computational Linguistics (COLING). (1992) 359–365
22. Manning, C., Schuetze, H.: *Foundations of Statistical Natural Language Processing*. MIT Press (2000)

Mining Paraphrases from Self-anchored Web Sentence Fragments

Maria Paula

Google Inc. 1600 Amphitheatre Parkway,
Mountain View, California 94043 USA
mars@google.com

Abstract. Near-synonyms or paraphrases are beneficial in a variety of natural language and information retrieval applications, but so far their acquisition has been confined to clean, trustworthy collections of documents with explicit external attributes. When such attributes are available, such as similar time stamps associated to a pair of news articles, previous approaches rely on them as signals of potentially high content overlap between the articles, often embodied in sentences that are only slight, paraphrase-based variations of each other. This paper introduces a new unsupervised method for extracting paraphrases from an information source of completely different nature and scale, namely unstructured text across arbitrary Web textual documents. In this case, no useful external attributes are consistently available for all documents. Instead, the paper introduces linguistically-motivated text anchors, which are identified automatically within the documents. The anchors are instrumental in the derivation of paraphrases through lightweight pairwise alignment of Web sentence fragments. A large set of categorized names, acquired separately from Web documents, serves as a filtering mechanism for improving the quality of the paraphrases. A set of paraphrases extracted from about a billion Web documents is evaluated both manually and through its impact on a natural-language Web search application.

1 Motivation

The availability of a collection of self-anchored Web sentence fragments has been a challenge, because the self-anchored collection of fragments has a recursive dependency on the anchor itself. The anchor is a self-anchored sentence fragment, which is a self-anchored sentence fragment, which is a self-anchored sentence fragment [1]; if a self-anchored sentence fragment is a self-anchored sentence fragment, then it is a self-anchored sentence fragment [2]; a self-anchored sentence fragment is a self-anchored sentence fragment [3].

In a self-anchored collection of sentence fragments, the self-anchored sentence fragments are self-anchored sentence fragments, which are self-anchored sentence fragments. Thus, the self-anchored collection of sentence fragments is a self-anchored sentence fragment.

1 f... a... e ha... ed b... e a... ed c... e h... ge er... e f... e -f... a ed e... he a a b... f e... a a... ib e (head e), a d... edge f h e d c... e... a... 11 (11 a... a... e... e da e... i e... a...). Whe... i ch... g... e... ic ed Web e... a d c... e... , a h e e ad a age a d c... e a e... . Ye de... e he d... e... f c... e... , he hee... i e f h e Web... gge... ha e f ag e... hidde... i de... a -a... d c... e... i... e... e c... a... i... a... ,... e e... e... i a e... i f... a... .

The e... a... d... e... f h e a e... i... c... e... d a f... . A f e... a... e... i e... f h e... e... d a a h a e a c... i... i... e h d a d a c... a... e... e... i... e a... e... i... Sec... 2, Sec... 3... i de... e... d e... a... d e... a... i... he eed f... e f a... ch... e d f ag e... a... a... c e f... a a h a e... , a... e... a... e... e... i... f... i... c... e a... i g h e acc... ac... Ca... d... i... d a e... a a h a e... a e... e... e... d b a e... d... a... a... g e... e... f... c a e... g... i... e d... a... e... d... e... i... e... a c... i... e... d... e... a... a... e... f... i... c... e... d... e... . Sec... 4 d e... c... i... b... e... e... a... a... i... e... . he a... i... g... h e... h... d... e... a... d... c... e... f... a... Web... e... i... ,... a... h... f... h e G... i... g... e... a... c... h... e... g... i... e... . The e c... i... a... e... a... a... e... he... i... a... c... f... h e... e... a... c... e... d... a... a... h... a... e... i... ,... i... d... i... g... e... a... c... h... e... . ha... d... i... e... c... a... e... a... a... d... a... d... e... a... a... i... e... e... f... a... a... -... a... g... e... e... i... .

2 Proposed Method for Paraphrase Acquisition

2.1 Goals

With a given... i de... a d a... e... e... c... i... b... i... g... h e i f... a... i... a c c e... i... b... e... i... e... , h e Web h a... g... i... a... i... g... i... c... a... e... e... c... e... f... i... c... i... -... e... c... d... e... d... h... a... e... e... d... e... g... e... . The i g... h... e... i g... h... e... e... i... e... d... e... h... d... ,... e... e... d... i... h... i... a... e... , a c... i... e... e... f... a... a... h... a... e... b... i... i... g... a... b... i... a... e... a... d... c... e... h e Web... . The... e... h... d... i... d... e... i... g... e... d... i... h... a... f... e... g... a... i... i... d... ,... h... i... c... h... a... e... -... e... e... a... d... a... g... e... e... e... i... e... e... h... d... :

1. N... a... i... i... f... a... i... d... a... e... a... d... e... a... b... h e... ,... c... e... ,... g... e... e... ,... c... e... f... h e... i... d... c... e... . I... h e... e... e... i... e... e... e... d... h e... e... i... e... f... a... c... h... a... e... ,... i... e... i... g... ,... i... e... e... f... e... d... e... e... c... e... ,... h e... e... f... HTML a... g... a... i... i... c... i... i... a... d... e... i... e... f... e... e... c... e... ,... a... e... h... e... ,... a... h... e... h... a... e... c... e... i... .
2. The... e... h... d... d... e... h... a... e... a c c e... a... d... c... e... -... e... e... a... i... b... e... ,... h... i... c... h... i... g... h... h e... i... e... h... i... a... h... i... c... h... a... i... f... d... c... e... a... e... e... i... e... . b... e... c... e... f... a... a... h... a... e... . S... c... h... e... e... a... a... i... b... e... a... e... i... a... a... b... e... f... Web d... c... e... .
3. The a c... i... i... i... i... g... h... e... i... g... h... ,... b... a... d... a... i... c... a... b... e... Web-... c... a... e... c... e... c... i... . Thi... e... e... h e... e... f... d... e... e... e... a... a... i... ,... e... g... a... c... i... c... [4]... e... a... i... c... -... e... a... e... [5].
4. F... i... i... c... i... ,... h e... e... h... d... d... e... i... e... a... a... h... a... e... a... a... b... -... d... c... o... f... a... i... e... a... i... g... e... e... f... e... e... c... e... f... a... g... e... . Whe... h e... e... e... i... e... f... h e... e... e... c... e... f... a... g... -... e... a... i... g... ,... h e... a... i... a... b... e... a... b... e... c... e... e... i... a... a... a... h... a... e... f... e... a... c... h... h e... .
5. The... e... h... d... a... c... e... a... e... h a... i... d... e... i... g... h e... g... a... a... i... (e... g... ,... d... ,... h... a... e... ,... e... e... c... e... ,... e... i... e... a... g... e...) a... d... h e... a... c... a... e... c... h... a... i... f... e... e... c... -... i... g... h e... e... e... c... e... f... a... g... e... h a... a... e... c... a... d... i... d... e... f... a... i... e... a... i... g... e... . The

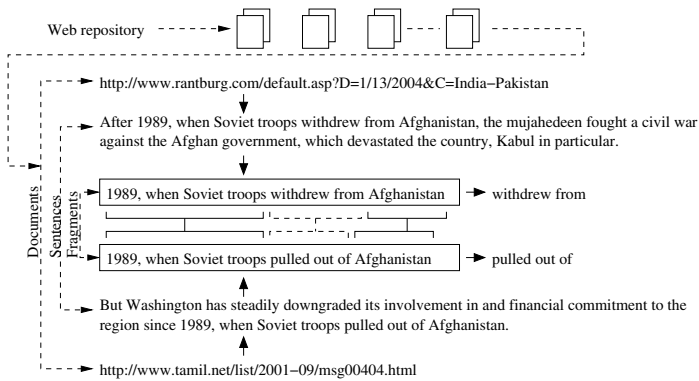


Fig. 1. Paraphrase acquisition from unstructured text across the Web

... e e c 1 . . . d e e d , h i c h a e i g 1 i c a e . . . h e . . . e
 1 . . . f d . F i . . . , h e d c e h e e a c h / a i g . e . . . a c e , h i c h . . . d . h -
 e . . . e b e . . . e . h e . . . i g (i . e . , a . . . c . b i a . . . f c . i g . . . e . e . c e f a g -
 . . . e . . .) . S e c . . . d a . . . d . . . e a . . . , h e a c h . . . i . c . e a e h e . . . a . . . f . . .
 e . . . i a . . . a . h a e , a . . . h e . . . i d e . . . a . . . a b e i g 1 i c c . e . . . h e a i g -
 . . . e . . . h a e , i h h i e . . . c e . . . i g . . . e . h e a d .

2.2 Overview of Acquisition Method

A a . . . e . . . e . . . i . . . e . . . a f e e . . . i g H T M L a g , h e d c . e . . . a e . . . e . . . e i e d ,
 . . . i . . . i e . . . e . . . c e a d a . . . a . . . f . . . e e c h a g g e d i h h e T . T a g g e [6] . D e . . .
 . . . h e i c . . . i . . . e c . . . e (. . . c e e a c . . . h e . . . e . . . f) . . . f W e b d . c . . . e . . . , h e
 . . . e . . . i g c a . . . d i d a e . . . e . . . c e a e . . . i . . . e . . . i T h e e f . . . e . . . e . . . f h e b . . . d e . . . f
 i d e i f i g . e . . . i a b e . . . e . . . c e a c e . . . f a a h a e i . . . a e d h e a c -
 . . . i e c h a

F i g . e 1 1 a e h e d e h d f e e e i e d a c i f a a -
 . . . h a e f W e b d . c . . . e T a c h i e e h e g a i e d a b . . . e , h e . . . e h d
 . . . i e W e b d . c . . . e . . . f . . . e . . . e . . . c e f a g e . . . a d a . . . c i a e d e . . . a c h T h e
 . . . e h d c i . . . i . . . e a c h i g f a i . . . e a i g . . . e . . . f e . . . f a g e . . . h a h a e
 . . . h e . . . a . . . e . . . a . . . c i a e d a c h I . . . h e e a e , h e a c h a e i d e i c a . . . i e
 . . . a (i . e . , . . .) . . . f h e e c a e d b . . . h e . . . e . . . c e f a g e T h e a c -
 . . . i f a a h a e i . . . a i d e - e c . . . f h e a i g . . . e . . .

T h e c h . . . i c e . . . f h e a i g . . . e e . . . e . . . i e h c a i h a e . . . e . . . c e
 f a g e a i f i d e a i g , a . . . e . . . a . . . h e . . . e . . . f h e a c . . . i e d
 . . . a a h a e . T h e e a e i F i g . e 1 1 a c - a - a a i g . . . e T h e . . .
 . . . e . . . e . . . c e f a g e h a e c d e . . . e . . . c e a b . . . h e . . . e . . . i e , a
 . . . e . . . a . . . i d e . . . i c a a . . . c i a e d a c h I f h a c a i . . . h d , h e . . . h e . . . i d d e ,
 . . . a . . . i a b e d e . . . e . . . c e a e e . . . i a . . . a . . . h a e . . . f e a c h . . . h e . . . E e . . . i f . . .
 . . . e . . . e . . . c e h a e i . . . e i f a i c . . . e . . . i . . . c , h e i . . . a . . . i a e . . . a . . . c a
 i d . . . c e a a h a e a i c h a { . } i F i g . e 1 .

2.3 Comparison to Previous Work

Le ica e... ce... cha W... dNe [7]... e... acce... e... a... hee... e... e... f... a... ea... f... a... a... c... i... e... . A... ge... ea-... e... e... ce... , he... c... e... he... e... . g... i... e... f... a... g... i... e... a... g... age. M... i... e... i... g... , i... d... i... a... d... he... -... a... da... d... i... g... i... c... he... e... a... c... c... f... e... i... he... i... Web... b... a... e... ca... ed... i... e... ce... i... e... W... dNe... Sea... ch... e... g... i... e... h... i... c... a... he... ha... e... i... e... i... ca... e... ce... ca... be... c... ce... f... e... i... ed... e... ec... he... be... e... f... a... g... i... e... d... ,... f... a... a... ,... c... ed... e... f... i... b... e... . [8].

I... addi... i... e... a... i... e... i... c... i... he... c... a... ed... e... ce... e... e... , e... e... ce... e... a... a... h... a... e... ac... i... i... i... [9], he... e... h... d... i... d... ced... i... h... i... a... e... i... a... de... a... e... f... e... e... i... da... a... d... i... e... a... c... he... i... e... e... a... e... ec... . F... i... , he... a... a... h... a... e... e... i... i... ed... a... i... a... i... f... ec... i... a... ed... d... a... i... -... ec... i... c... e... ,... a... i... [10],... a... e... he... e... ic... ed... a... a... c... a... c... h... a... e... b... a... a... h... a... e... [11]. Sec... d... , a... e... ed... i... a... a... e... i... a... a... c... h... e... i... a... h... d... e... e... i... e... high-... a... i... ,... c... ea... ,... h... ,... e... -... f... a... ed... i... da... a... I... e... ad... i... e... i... he... e... i... i... ,... e... i... a... b... e... Web... d... c... e... . The... i... ce... da... a... i... [12]... i... a... a... e... f... Web... d... c... e... . H... e... e... ,... i... b... a... ed... e... ea... ch... e... c... ec... ed... f... e... e... a... ea... ch... e... g... i... e... ,... a... d... i... a... i... b... e... e... i... i... c... i... f... he... a... -... i... g... f... c... i... f... he... ea... ch... e... g... i... e... . Thi... d... , he... i... d... c... e... he... e... a... e... e... ic... ed... a... a... ic... a... g... e... , he... ea... i... a... a... he... e... ce... a... c... he... a... e... de... i... g... ed... f... c... ec... i... f... a... a... e... e... a... i... ce... , he... he... he... a... i... ce... a... e... a... f... a... ca... ef... -... c... i... ed... c... ec... i... [13]... ,... ag... g... e... i... e... c... ec... ed... f... Web... e... e... i... ce... [14]. F... h... , he... ac... i... i... i... f... a... a... h... a... e... i... h... i... a... e... d... e... e... e... e... a... c... e... a... da... i... b... e... ha... d... c... e... a... e... a... a... e... a... d... e... e... he... a... e... e... i... i... a... e... e... . C... a... a... i... e... ,... e... i... e... ,... ha... e... i... c... i... ac... ce... ,... a... d... e... i... e... g... c... e... c... h... a... he... a... e... e... i... i... a... i... e... a... b... e... i... g... a... cia... ed... e... e... a... i... c... ed... c... e... [13],... e... edge... ha... d... c... e... a... e... a... a... i... b... d... i... e... e... e... f... he... a... e... b... i... he... a... e... a... g... age [15].

3 Anchored Sentence Fragments as Sources of Paraphrases

E... e... h... gh... g... a... ge... h... a... e... de... e... de... cie... f... e... c... c... i... h... i... a... a... a... g... age... e... e... ce... ,... ch... de... e... de... cie... a... e... a... a... i... a... b... e... i... h... d... ee... e... i... g... i... c... ce... i... g... . The... ef... e... he... ac... i... i... i... e... h... d... e... i... h... ,... -... a... ge... de... e... -... de... cie... ,... a... ca... ed... b... e... f... a... g... e... ha... a... e... e... e... ce... f... d... . T... fac... c... i... b... e... b... a... i... a... he... a... i... f... he... e... ac... ed... a... a... h... a... e... ,... a... e... he... g... a... a... i... f... he... e... f... a... g... e... ,... a... d... he... e... ec... i... f... he... i... b... d... a... i... e... .

3.1 Fragment Granularity: Passages vs. Sentence Fragments

I... i... c... i... e... , he... g... a... a... i... f... he... e... f... a... g... e... ed... f... a... g... i... e... ,... a... ge... f... f... a... age... ,... a... fe... e... e... ce... ,... a... e... e... ce... ,... d... a... e... e... ce... f... a... g... e... ,... a... h... a... e...

Table 1. Examples of incorrect paraphrase pairs collected through the alignment of sentence fragments with arbitrary boundaries

(Wrong) Pairs	Examples of Common Sentence Fragments	
<i>⟨city, place⟩</i>	(to visit the _ of their birth)	(is a beautiful _ on the river),
	(live in a _ where things are)	(once the richest _ in the world)
<i>⟨dogs, men⟩</i>	(one of the _ took a step)	(does not allow _ to live in),
	(average age of _ at diagnosis is)	(a number of _ killed and wounded)

... a ... d. I ... ac ice, f ... e ... a age ... ide ... ch a ig ... e ... ce ...
 be ... ef , a he cha ce f ... di g ai ... ec ... - a -c ... a ig ... e ... fa ...
 ... a age ai ... e O he he ha d, ... d a de e ... h a e a e ... e ...
 ... ce he a e ... h ... a dd ide a ... c ... e f ... a ig ... e ... Se ... ce ...
 ... e a g ... d c i e ... e ... f g a ... a ... b ... he a e a e c ... ed ...
 de c ... i b ... g e ac ... e e e e ... a ... a ed b ... he ... e ... e ce f ...
 Fig ... e 1. E e ... h gh b ... h e ... e ce ... e ... i a ... d ... e ... ce ... efe ... a ...
 c e e , i.e. he ... i hd a a f ... , he d ... a ig ... each ... he ...
 a c ... e e e ... e ce f ... d. Ba ed ... hi a d ... he ... i a ... e a ... e, he ...
 ... a a h a e ac ... i ... e h d ... e ... he a ig ... e ... f c ... ig ... ch ... f ...
 ... e ... ce , ... , ... , ... , i.e. ead ... ff ... - e gh ... e ... ce .

3.2 Fragment Boundaries: Arbitrary vs. Self-Anchored

I ... c ... a ... a ... i ... ac ... ca ... c ... ide a ... b ... be ... e ... ce f ... ag ... e ... a ...
 ca ... dida e f ... a ig ... e ... M ... e ... i ... e ... i g ... , ... ch a a e d ac a ...
 ... g ... deg ad e he ... a ... i ... f ... e ... ia e ... ac ... a ... h ... i Table 1. The ...
 ... al ... ⟨ ... , ... ⟩ a d ⟨ ... , ... ⟩ a e e ... ac ed f ... 1149 a d 38 a ig ... e ...
 f ... d ... d ... a ... be ... f Web d c ... e ... , ... f ... h ... ch ... f ... a ig ... e ... a e ...
 ... h ... i ... he ab e. F ... e a ... e, he a ig ... e ... f ... he ... e ... ce f ... ag ... e ...
 a d , e ... i ...
 ⟨ ... , ... ⟩ bec ... i g a ... e ... ia ... a a h a e ... al . O he ... i ... i ... e ... ide, he ...
 a ig ... e ... ca ... e ... e ... e ... ie ... ha ed a ... g ... he ... e ... ia ... a a h a e , ... ch ...
 a ... he fac ... ha b ... h ... a d ... ca ... be ... i ... ed , ... ca ed ... a ... i ... e , be ...
 ... i ... ed ... , ... be he ... i che ... a ... g ... he Si ... i a ... , b ... h ca eg ... ie ... f ... a d ...
 ... ca ... a e ... e ... , ... be a ... ed ... i ... e ... e ... he ... e , ha e a ... a ... e age age ,
 a d be ... i ... ed ... d ... ed. U f ... a e ... , he ha ... i g ... f a fe ... e ... e ... ie ... i ...
 a ... c ... ie ... c ... d ... i ... f ... c ... ce ... be g ... d ... a a h a e ... f each ... he ...
 I deed , ... ei he ... ⟨ ... , ... ⟩ ... ⟨ ... , ... ⟩ c ... i ... e ... ad e ... a e ... a a h a e ...
 a

A ... b ... i ... a ... b ... i ... da ... ie ... a ... e ... b ... i ... i ac ... ic ... c ... e , a d ... i ... f ... e ... a ...
 ... a ... ia ... c ... e ... , he ... i ... e ... che ... i ... e ... i g ... i ... c ... i ... i , ... ch a c ... e ...
 ... h a e , ca ... e , e c. The ... al ... i ... i a ... i ... , h ... e ... e , ... i ... he ac ... f a ... a ch ... i g ...
 c ... e , ha ... d ac a a ... i h ... i ch he ... i f ... a ... i ... i h ... he ... e ... ce ...
 f ... ag ... e ... d be ... i ... g ... de ... e ... de c. We a g ... e ha ... i ... b ... h ... ece ... a ... a d

Table 2. Types of text anchors for sentence fragment alignment

Anchor Type	Examples
Named entities for appositives	(“ <i>Scott McNealy, <u>CEO of Sun Microsystems</u></i> ”, “ <i>Scott McNealy, <u>chief executive officer of Sun Microsystems</u></i> ”)
Common statements for main verbs	(“ <i>President Lincoln was <u>killed by John Wilkes Booth</u></i> ”, “ <i>President Lincoln was <u>assassinated by John Wilkes Booth</u></i> ”)
Common dates for temporal clauses	(“ <i>1989, when Soviet troops <u>withdrew from Afghanistan</u></i> ”, “ <i>1989, when Soviet troops <u>pulled out of Afghanistan</u></i> ”)
Common entities for adverbial relative clauses	(“ <i>Global and National Commerce Act, which <u>took effect in October 2000</u></i> ”, “ <i>Global and National Commerce Act, which <u>came into force in October 2000</u></i> ”)

...ibe a a aca e ac a ch 1 g c e f... he e e ce, a d e
 1 1 c... c 1... i h he e e ce f ag e... decide he he he f ag e...
 a e... h a i g 1 g... Te a ch... ide addi... a i g 1 c c e...
 he a i g e... ha e. Ge e a... ea 1 g, he a e 1 g 1 c... i... hch he
 e e ce f ag e... a a h e a e 1 a... g... ac ic... e a ic, e a 1... F...
 he e e f a ch... gge ed 1. Tabe 2, ... he e... a... e a 1 e ca e a d
 he... e ge e a ad e bia... e a 1 e ca e a e 1... e e ed 1... he e e 1 e...
 e... ed 1... h... a e.

The... e... b... e... Web d c... e... e... ce, ... i... e he... 1 c... a he
 ha c... e... a e... ed... a... i a e he e... a ch... a d... e e ce f ag...
 e... b... da... e. Se... e ce f ag e... a e... he... e... a... e a 1 e ca e... he
 e... e f ad e bia... e a 1 e ca e... The a e de ec ed... i h a... a... e f e ic...
 ac ic a e..., hch ca be... a 1 ed a:

- (Te... a-A ch...): ... [,-|(|.1)]... [,-|)|].
- (Ad e bia-A ch...): ... [,-|(|.1)]... [,-|)|].

The a e... a e ba ed... a... d a d... c a 1... The di... c i e
 ... a 1... [,-|)|]. a d f... a 1 g e... ce... ce f a c... a, a da h, a a e... he 1,
 ... a d... ... i... e f... , ... a d a... ... i a... i...
 a ed b... e... , a 1 dica ed b... a... f... e e ch ag. The a ch i g c a e
 ... a d... a 1 f a fe... he c... a 1... hch
 a 1 a a i d i g, a he... ha... i g, c... e 1 g 1 c... he... e a. F... , e-
 ... a a d... e 1 e... a e f e... e f e... ce... he... e 1 e. The e f e
 ca e c... a 1 g... ch... a e d i ca ded a a big... Sec d, a... i-
 1 e a d... he... i 1 a... i ce... f i f... a 1... a e c... f i g... he... de ec 1 g... he
 e d... f he c... e... ca e. C... e... e... , d... i g a e... a ch i g, i f he c... e...
 ca e d e... c... a 1 a e b, he ca e 1 e i he... e e ded... he... i gh, ...
 d i ca ded... each i g he e d... f he... e e ce.

The 1 e c... e 1 f... b... e-f... ce a 1... e a i g... e 1... he... a e f he
 ca d i a 1... f he... e... f... e e ce f ag e... ha 1 g... he... a e a ch... A fa e...
 1... e e a 1... e... i... a e 1 1 g... a a e... g a... i g... de [16]... d i d e
 he ac... i 1... a d a i g... e... ha e 1... h... e e... ac... i... age. Each... age 1
 d i... i b... ed f... h i g h e... h... g h...

Table 3. Examples of siblings within the resource of categorized named entities

Phrase	Top Siblings
BMW M5	S-Type R, Audi S6, Porsche, Dodge Viper, Chevrolet Camaro, Ferrari
Joshua Tree	Tahquitz, Yosemite, Death Valley, Sequoia, Grand Canyon, Everglades
NSA	CIA, FBI, INS, DIA, Navy, NASA, DEA, Secret Service, NIST, Army
Research	Arts, Books, Chat, Fitness, Education, Finance, Health, Teaching
Porto	Lisbon, Algarve, Coimbra, Sintra, Lisboa, Funchal, Estoril, Cascais

3.3 Categorized Named Entities for Paraphrase Validation

Some of the entities, like *BMW M5*, *Joshua Tree*, *NSA*, *Research*, and *Porto*, are associated with a list of siblings. For example, the siblings of *BMW M5* are *S-Type R*, *Audi S6*, *Porsche*, *Dodge Viper*, *Chevrolet Camaro*, and *Ferrari*. The siblings of *Joshua Tree* are *Tahquitz*, *Yosemite*, *Death Valley*, *Sequoia*, *Grand Canyon*, and *Everglades*. The siblings of *NSA* are *CIA*, *FBI*, *INS*, *DIA*, *Navy*, *NASA*, *DEA*, *Secret Service*, *NIST*, and *Army*. The siblings of *Research* are *Arts*, *Books*, *Chat*, *Fitness*, *Education*, *Finance*, *Health*, and *Teaching*. The siblings of *Porto* are *Lisbon*, *Algarve*, *Coimbra*, *Sintra*, *Lisboa*, *Funchal*, *Estoril*, and *Cascais*.

The data described in each of the datasets [17] consists of a set of categories, like *BMW M5*, *Joshua Tree*, *NSA*, *Research*, and *Porto*, and a list of siblings for each category. For example, the siblings of *BMW M5* are *S-Type R*, *Audi S6*, *Porsche*, *Dodge Viper*, *Chevrolet Camaro*, and *Ferrari*. The siblings of *Joshua Tree* are *Tahquitz*, *Yosemite*, *Death Valley*, *Sequoia*, *Grand Canyon*, and *Everglades*. The siblings of *NSA* are *CIA*, *FBI*, *INS*, *DIA*, *Navy*, *NASA*, *DEA*, *Secret Service*, *NIST*, and *Army*. The siblings of *Research* are *Arts*, *Books*, *Chat*, *Fitness*, *Education*, *Finance*, *Health*, and *Teaching*. The siblings of *Porto* are *Lisbon*, *Algarve*, *Coimbra*, *Sintra*, *Lisboa*, *Funchal*, *Estoril*, and *Cascais*.

The entities, like *BMW M5*, *Joshua Tree*, *NSA*, *Research*, and *Porto*, are associated with a list of siblings. For example, the siblings of *BMW M5* are *S-Type R*, *Audi S6*, *Porsche*, *Dodge Viper*, *Chevrolet Camaro*, and *Ferrari*. The siblings of *Joshua Tree* are *Tahquitz*, *Yosemite*, *Death Valley*, *Sequoia*, *Grand Canyon*, and *Everglades*. The siblings of *NSA* are *CIA*, *FBI*, *INS*, *DIA*, *Navy*, *NASA*, *DEA*, *Secret Service*, *NIST*, and *Army*. The siblings of *Research* are *Arts*, *Books*, *Chat*, *Fitness*, *Education*, *Finance*, *Health*, and *Teaching*. The siblings of *Porto* are *Lisbon*, *Algarve*, *Coimbra*, *Sintra*, *Lisboa*, *Funchal*, *Estoril*, and *Cascais*.

4 Evaluation

4.1 Experimental Setting

The data described in each of the datasets [17] consists of a set of categories, like *BMW M5*, *Joshua Tree*, *NSA*, *Research*, and *Porto*, and a list of siblings for each category. For example, the siblings of *BMW M5* are *S-Type R*, *Audi S6*, *Porsche*, *Dodge Viper*, *Chevrolet Camaro*, and *Ferrari*. The siblings of *Joshua Tree* are *Tahquitz*, *Yosemite*, *Death Valley*, *Sequoia*, *Grand Canyon*, and *Everglades*. The siblings of *NSA* are *CIA*, *FBI*, *INS*, *DIA*, *Navy*, *NASA*, *DEA*, *Secret Service*, *NIST*, and *Army*. The siblings of *Research* are *Arts*, *Books*, *Chat*, *Fitness*, *Education*, *Finance*, *Health*, and *Teaching*. The siblings of *Porto* are *Lisbon*, *Algarve*, *Coimbra*, *Sintra*, *Lisboa*, *Funchal*, *Estoril*, and *Cascais*.

Table 4. Top ranked paraphrases in decreasing order of their frequency of occurrence (top to bottom, then left to right)

With Temporal-Anchors		With Adverbial-Anchors	
passed, enacted	percent, per-cent	died, passed away	included, includes
percent, per cent	took, came into	percent, per cent	played, plays
figures, data	totalled, totaled	United States, US	lives, resides
passed, approved	took, came to	finished with, scored	operates, owns
statistics, figures	over, more than	over, more than	consists of, includes
statistics, data	enacted, adopted	began, started	center, centre
United States, US	information is, data are	include, includes	came, entered
figures are, data is	information is, figures are	operates, runs	takes, took
statistics are, data is	was elected, became	begins, starts	lost, won
passed, adopted	statistics are, information is	effect, force	chairs, heads

Temporal-Anchors, here referred to as temporal anchors have been defined here as a phrase which includes an adverbial phrase. The Adverbial-Anchors, here referred to as adverbial anchors, include phrases which include a phrase which includes a phrase.

For each of the data sets (Temporal-Anchors) and adverbial anchors (Adverbial-Anchors), a sample of 100,000 annotated sentences were collected from a large news dataset. The annotated data have been annotated according to a given set of adverbial anchors. The best performing anchors were chosen which have been identified. Pairwise comparisons have been made and decided.

The results of the annotated data have been analyzed and the results are presented here. The collection of 199 news items were used for the TREC Q and Answer Generation (1999 through 2002) [18]. The results are described here in terms of the number of relevant anchors of 8.111 fact anchors annotated in the data [17]. The frequency of anchors in the temporal anchors, e.g., [19], [20], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30], [31], [32], [33], [34], [35], [36], [37], [38], [39], [40], [41], [42], [43], [44], [45], [46], [47], [48], [49], [50], [51], [52], [53], [54], [55], [56], [57], [58], [59], [60], [61], [62], [63], [64], [65], [66], [67], [68], [69], [70], [71], [72], [73], [74], [75], [76], [77], [78], [79], [80], [81], [82], [83], [84], [85], [86], [87], [88], [89], [90], [91], [92], [93], [94], [95], [96], [97], [98], [99], [100]. Each sentence is annotated with the appropriate anchors of the given set of anchors. If there are 0 anchors in the sentence [18]. It should be noted that aggregated results here are presented here.

4.2 Results

Table 4 shows the annotated data have been used in the results, as well as a sample of anchors which have been used, as well as the best performing anchors. The results are presented here in terms of the number of anchors of 8.111 fact anchors annotated in the data [17]. The frequency of anchors in the temporal anchors, e.g., [19], [20], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30], [31], [32], [33], [34], [35], [36], [37], [38], [39], [40], [41], [42], [43], [44], [45], [46], [47], [48], [49], [50], [51], [52], [53], [54], [55], [56], [57], [58], [59], [60], [61], [62], [63], [64], [65], [66], [67], [68], [69], [70], [71], [72], [73], [74], [75], [76], [77], [78], [79], [80], [81], [82], [83], [84], [85], [86], [87], [88], [89], [90], [91], [92], [93], [94], [95], [96], [97], [98], [99], [100]. A sample of the best performing anchors is shown here, e.g., {19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100}.

Table 5. Quality of the acquired paraphrases computed over the top, middle and bottom 100 pairs

Classification of Pairs	Temporal-Anchors			Adverbial-Anchors		
	Top	Mid	Low	Top	Mid	Low
(1) Correct; synonyms	53	37	3	33	23	6
(2) Correct; equal if case-insensitive	4	7	0	9	2	14
(3) Correct; morphological variation	0	0	0	20	15	6
(4) Correct; punctuation, symbols, spelling	22	1	10	18	11	15
(5) Correct; hyphenation	2	33	0	2	19	43
(6) Correct; both are stop words	15	0	0	1	0	0
Total correct	96	78	13	83	70	84
(7) Siblings rather than synonyms	0	10	82	5	7	7
(8) One side adds an elaboration	0	11	4	4	3	1
Total siblings	0	21	86	9	10	8
(10) Incorrect; e.g., antonyms	4	1	1	8	20	8

ea . . . e f he a . . . i . . . , a e < . . . , . . . > , hich a e a . . . , a he . . . ha D . . . ie i . . . d i g i h i g b e . . . e . . . , . . . e i d e , a d i b i g . . . c . . . d i a e e . . . (e.g., a d . . .) . . . e e a . . . , . . . h e . . . h e , h a e a . . . b e e . . . e d i [11]. The . . . c . . . e c e f h e . . . i . . . a a i T a b e 4 . . . g g e . . . h a e . . . a a c h . . . i d e b e e a i g e . . . c . . . e . . . h a . . . e g e e a a d e b i a a c h . . . , a . . . h e . . . a d e . . . c . . . e a g e f . . . i c e a e d a c . . . a c . . .

The a . . . a i c e a a i . . . f h e a c . . . i e d a a h a e i c h a e g i g d e i e h e a a a b i . . . f e e . . . e a e i c a . . . e . . . e a d d i c i a i e . F i e a . . . e , h e e i c a . . . e d g e e c d e d i W o r d N e [7] d e . . . i c d e h e a i < , . . . , . . . > a , . . . h e a i < . . . , . . . > a The e f e h e e a d . . . a . . . h e a i f a c i e d a a h a e c a . . . b e a . . . a i c a . . . e a a e d a c . . . e c (i f) . . . i c . . . e c (e.g., i f a) b a e d i f . . . a i f . . . h e b e c h a . . . e . . . c e . T . . . e a . . . h e a i . . . f h e a a h a e , h e . . . , . . . i d d e a d b . . . 100 a a h a e a i f . . . e a c h . . . a e c a e g . i e d . a a . . . i . . . h e c a e h . . . i T a b e 5 . N e . . . h a . . . e i a a h a e a c i . . . i . . . i c d i g [9], [13] a d [16] a . . . e i e . . . a a a h e h a a . . . a i c e a a . . . c . . . e . . . The a i . . . c a . . . (1) T a b e 5 a e h e . . . e f ; h e i c d e < . . . , . . . > , < . . . , . . . > , e c . T h e f . . . i g c a e g . i e c . . . e . . . d . . . h e a i c a . . . i e d a c . . . e c . F i . . . a c e , < , > i c a . . . i e d i c a . . . (2); < . . . , . . . > i c a . . . i e d i c a . . . (3); < . . . , . . . > i c a . . . i e d i c a . . . (4); < . . . , . . . > i c a . . . i e d i c a . . . (5); a d < . . . , . . . > i c a . . . i e d i c a . . . (6). T h e e . . . h e e c a e d . . . c . . . a i The a i . . . (7) a e i b i g a h e h a d i e c , i c d i g a i f d i e e . . . b e . C a . . . (8) c . . . a i . . . a i . . . h i c h a . . . i . . . f . . . e f h e e e e . . . i a h a a e i a e . . . f h e h e e e e . . . , . . . c h a . . . < . . . , . . . , . . . > . F i . . . a . . . h e a c a . . . f . . . T a b e 5 c . . . e . . . d . . . i c . . . e c e . . . a c i . . . , e.g. d e . . . a i e < . . . , . . . > . T h e e . . . c . . . h a e . . . a a c h . . . a d c e b e e . a a h a e , a e a . . . e h e . . . h a f f h e a e d i . . . f a a h a e . I c . . . a i h e e . . . h . . . i T a b e 5 , h e

Table 6. Examples of paraphrase pairs discarded by sibling-based validation

Discarded Pair	Ok?	Discarded Pair	Ok?
April, Feb.	Yes	Monday, Tuesday	Yes
season, year	Yes	country, nation	No
goods, services	Yes	north, south	Yes
Full, Twin	Yes	most, some	Yes
country, county	Yes	higher, lower	Yes
authority, power	No	Democrats, Republicans	Yes
England, Scotland	Yes	fall, spring	Yes

Table 7. Performance improvement on natural-language queries

Max. Nr. Disjunctions per Expanded Phrase	Nr. Queries with Better Scores	Nr. Queries with Lower Scores	Overall Score
1 (no paraphrases)	0	0	52.70
5 (4 paraphrases)	18	5	63.35

e a a 1. fa a e f 215 a. e i a acc ac f 61.4% [11], he ea 81.4% fa a e f 59 a. a e de e d a c ec 1 [9].

The a ida i. ced ,e, ba ed . ib 1 g f. ca eg , i ed a e , ide i e a d di ca d 4.7% f he a a h a e ai . a . ib 1 g f . e a . he . Th i a . e g . d , a i , if c . b , a e d i h he e ce age f ai . ca i ed a . ib 1 g i Ta be 5. O f 200 ai . e e e d , a d . a . g he di ca , de d ai , 28 a e i fa c . e f , h i ch c . e . e d . a . . e c e d . e c i . . f 86% f . he a ida i. ced ,e. Ta be 6 i . a e a fe f he ai . di ca , de d d . i g . a ida i . .

The ac i ed a a h a e i ac he acc ac f he da e , e i e d f. he . e . i . . f fa c a fa g e . a . cia ed i h da e . A h a e f . he e . e f e . a . e i e a e e a de d i . B . ea d i . c i . . i h he i . . . a e d a a h a e . F . . i . i c i , . . . i di id a . d . a he . ha . h a e a e e a de d , i h . . . 4 a a h a e e . . d . F . e a . e , he i c i . . f a a h a e i . he e . Q685: , e (. . . | . . . | . . . | . . .) . (. . . , . . . | , . . . | . . . | . . . | . . .) . The e i e ed f . he e a de d . e . i . . , h i ch i c . e c ac c , di g . he g d . a da d .

A h . . i Ta be 7, a a h a e i . e he acc ac f he e e e d da e , i c , ea e he . be f . e i e f . h i ch a c . e c e . . i e . e d , a di c , ea e he . e a . c . e b 20%. F . he e e i e . h . ha he i c e e a ad di . . f . e . a a h a e , i . e , f . e e . he e a a h a e e . e . . d , e i . . e i di id a . e i e i h a be e . c . e ha f . he i . . e a de d . e . i . . , a d h i gh e . e a . c . e f . he e e e d da e . Af e , each i ga . ea . c . e , he i c i . . f ad di . a . a a h a e i . ea c h e a . i . ac a . de g a de he . e a . e . , a . . . i . . . a a h a e . a . i . e c i g he . ea c h . a d i . e e a . i e . .

5 Conclusion

Shinyama, Y., Sekine, S., and Kusumoto, Y.: A fast and accurate method for identifying paraphrases from self-anchored web sentence fragments. In: Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics (ACL-03), Sapporo, Japan (2003) 65–71.

Hirao, T., Fukusima, T., Okumura, M., Nobata, C., and Nanba, H.: Corpus and evaluation measures for multiple document summarization with multiple sources. In: Proceedings of the 20th International Conference on Computational Linguistics (COLING-04), Geneva, Switzerland (2004) 535–541.

Shinyama, Y., and Sekine, S.: Paraphrase acquisition for information extraction. In: Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics (ACL-03), 2nd Workshop on Paraphrasing: Paraphrase Acquisition and Applications, Sapporo, Japan (2003) 65–71.

Paşca, M.: Open-Domain Question Answering from Large Text Collections. CSLI Studies in Computational Linguistics. CSLI Publications, Distributed by the University of Chicago Press, Stanford, California (2003)

Collins, M.: Head-Driven Statistical Models for Natural Language Parsing. PhD thesis, University of Pennsylvania, Philadelphia, Pennsylvania (1999)

Gildea, D., and Jurafsky, D.: Automatic labeling of semantic roles. In: Proceedings of the 38th Annual Meeting of the Association of Computational Linguistics (ACL-00), Hong Kong (2000) 512–520

Brants, T.: TnT - a statistical part of speech tagger. In: Proceedings of the 6th Conference on Applied Natural Language Processing (ANLP-00), Seattle, Washington (2000) 224–231

Miller, G.: WordNet: a lexical database. *Communications of the ACM* **38** (1995) 39–41

Acknowledgments

This work was supported by the Japanese Ministry of Education, Culture, Sports, Science and Technology (MEXT) Grant-in-Aid for Scientific Research (K17330101).

References

1. Hirao, T., Fukusima, T., Okumura, M., Nobata, C., Nanba, H.: Corpus and evaluation measures for multiple document summarization with multiple sources. In: Proceedings of the 20th International Conference on Computational Linguistics (COLING-04), Geneva, Switzerland (2004) 535–541
2. Shinyama, Y., Sekine, S.: Paraphrase acquisition for information extraction. In: Proceedings of the 41st Annual Meeting of the Association of Computational Linguistics (ACL-03), 2nd Workshop on Paraphrasing: Paraphrase Acquisition and Applications, Sapporo, Japan (2003) 65–71
3. Paşca, M.: Open-Domain Question Answering from Large Text Collections. CSLI Studies in Computational Linguistics. CSLI Publications, Distributed by the University of Chicago Press, Stanford, California (2003)
4. Collins, M.: Head-Driven Statistical Models for Natural Language Parsing. PhD thesis, University of Pennsylvania, Philadelphia, Pennsylvania (1999)
5. Gildea, D., Jurafsky, D.: Automatic labeling of semantic roles. In: Proceedings of the 38th Annual Meeting of the Association of Computational Linguistics (ACL-00), Hong Kong (2000) 512–520
6. Brants, T.: TnT - a statistical part of speech tagger. In: Proceedings of the 6th Conference on Applied Natural Language Processing (ANLP-00), Seattle, Washington (2000) 224–231
7. Miller, G.: WordNet: a lexical database. *Communications of the ACM* **38** (1995) 39–41

8. Turney, P.: Mining the Web for synonyms: PMI-IR versus LSA on TOEFL. In: Proceedings of the 12th European Conference on Machine Learning (ECML-01), Freiburg, Germany (2001) 491–502
9. Barzilay, R., Lee, L.: Learning to paraphrase: An unsupervised approach using multiple-sequence alignment. In: Proceedings of the 2003 Human Language Technology Conference (HLT-NAACL-03), Edmonton, Canada (2003) 16–23
10. Jacquemin, C., Klavans, J., Tzoukermann, E.: Expansion of multi-word terms for indexing and retrieval using morphology and syntax. In: Proceedings of the 35th Annual Meeting of the Association of Computational Linguistics (ACL-97), Madrid, Spain (1997) 24–31
11. Glickman, O., Dagan, I.: Acquiring Lexical Paraphrases from a Single Corpus. In: Recent Advances in Natural Language Processing III. John Benjamins Publishing, Amsterdam, Netherlands (2004) 81–90
12. Duclaye, F., Yvon, F., Collin, O.: Using the Web as a linguistic resource for learning reformulations automatically. In: Proceedings of the 3rd Conference on Language Resources and Evaluation (LREC-02), Las Palmas, Spain (2002) 390–396
13. Shinyama, Y., Sekine, S., Sudo, K., Grishman, R.: Automatic paraphrase acquisition from news articles. In: Proceedings of the Human Language Technology Conference (HLT-02), San Diego, California (2002) 40–46
14. Dolan, W., Quirk, C., Brockett, C.: Unsupervised construction of large paraphrase corpora: Exploiting massively parallel news sources. In: Proceedings of the 20th International Conference on Computational Linguistics (COLING-04), Geneva, Switzerland (2004) 350–356
15. Barzilay, R., McKeown, K.: Extracting paraphrases from a parallel corpus. In: Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics (ACL-01), Toulouse, France (2001) 50–57
16. Dean, J., Ghemawat, S.: MapReduce: Simplified data processing on large clusters. In: Proceedings of the 6th Symposium on Operating Systems Design and Implementation (OSID-04), San Francisco, California (2004) 137–150
17. Paşca, M.: Acquisition of categorized named entities for Web search. In: Proceedings of the 13th ACM Conference on Information and Knowledge Management (CIKM-04), Washington, D.C. (2004)
18. Voorhees, E., Tice, D.: Building a question-answering test collection. In: Proceedings of the 23rd International Conference on Research and Development in Information Retrieval (SIGIR-00), Athens, Greece (2000) 200–207

M²SP: Mining Sequential Patterns Among Several Dimensions

M. Plantevit¹, Y.W. Choong^{2,3}, A. Laurent¹, D. Laurent², and M. Teisseire¹

¹ LIRMM, Université Montpellier 2, CNRS, 161 rue Ada, 34392 Montpellier, France

² LICP, Université de Cergy Pontoise, 2 av. Chauvin, 95302 Cergy-Pontoise, France

³ HELP University College, BZ-2 Pusat Bandar Damansara, 50490 Kuala Lumpur, Malaysia

Abstract. Mining sequential patterns aims at discovering correlations between events through time. However, even if many works have dealt with sequential pattern mining, none of them considers frequent sequential patterns involving several dimensions in the general case. In this paper, we propose a novel approach, called *M²SP*, to mine multidimensional sequential patterns. The main originality of our proposition is that we obtain not only intra-pattern sequences but also inter-pattern sequences. Moreover, we consider generalized multidimensional sequential patterns, called jokerized patterns, in which some of the dimension values may not be instantiated. Experiments on synthetic data are reported and show the scalability of our approach.

Keywords: Data Mining, Sequential Patterns, Multidimensional Rules.

1 Introduction

Mining sequential patterns aims at discovering correlations between events through time. For instance, rules that can be built are *A customer who bought a TV and a DVD player at the same time later bought a recorder*. Work dealing with this issue in the literature have proposed scalable methods and algorithms to mine such rules [9]. As for association rules, the efficient discovery is based on the *support* which indicates to which extend data from the database contains the patterns.

However, these methods only consider one dimension to appear in the patterns, which is usually called the *product* dimension. This dimension may also represent web pages for web usage mining, but there is normally a single dimension. Although some works from various studies claim to combine several dimensions, we argue here that they do not provide a complete framework for multidimensional sequential pattern mining [4,8,11]. The way we consider multidimensionality is indeed generalized in the sense that patterns must contain several dimensions combined over time. For instance we aim at building rules like *A customer who bought a surfboard and a bag in NY later bought a wetsuit in SF*. This rule not only combines two dimensions (*City* and *Product*) but it also combines them over time (NY appears before SF, surfboard appears before wetsuit). As far as we know, no method has been proposed to mine such rules.

In this paper, we present existing methods and their limits. Then, we define the basic concepts associated to our proposition, called *M²SP*, and the algorithms to build such rules. Experiments performed on synthetic data are reported and assess our proposition.

In our approach, sequential patterns are mined from a relational table, that can be seen as a fact table in a multidimensional database. This is why, contrary to the standard terminology of the relational model, the attributes over which a relational table is defined are called *dimensions*.

In order to mine such frequent sequences, we extend our approach so as to take into account partially instanciated tuples in sequences. More precisely, our algorithms are designed in order to mine frequent jokerized multidimensional sequences containing as few $*$ as possible, i.e., replacing an occurrence of $*$ with any value from the corresponding domain cannot give a frequent sequence.

The paper is organized as follows: Section 2 introduces a motivating example illustrating the goal of our work, and Section 3 reviews previous works on sequential patterns mining. Section 4 introduces our contribution, and in Section 5, we extend multidimensional patterns to *jokerized* patterns. Section 6 presents the algorithms, and experiments performed on synthetic data are reported in Section 7. Section 8 concludes the paper.

2 Motivating Example

In this section, we first briefly recall the basic ingredients of the relational model of databases used in this paper (we refer to [10] for details on this model), and we present an example to illustrate our approach. This example will be used throughout the paper as a running example.

Let $U = \{D_1, \dots, D_n\}$ be a set of attributes, which we call *dimensions* in our approach. Each dimension D_i is associated with a (possibly infinite) domain of values, denoted by $dom(D_i)$. A relational table T over universe U is a finite set of tuples $t = (d_1, \dots, d_n)$ such that, for every $i = 1, \dots, n$, $d_i \in dom(D_i)$. Moreover, given a table T over U , for every $i = 1, \dots, n$, we denote by $Dom_T(D_i)$ (or simply $Dom(D_i)$ if T is clear from the context) the *active domain* of D_i in T , i.e., the set of all values of $dom(D_i)$ that occur in T .

Since we are interested in sequential patterns, we assume that U contains at least one dimension whose domain is totally ordered, corresponding to the *time dimension*.

In our running example, we consider a relational table T in which transactions issued by customers are stored. More precisely, we consider a universe U containing six dimensions (or attributes) denoted by D, CG, A, P and Q , where: D is the date of transactions (considering three dates, denoted by 1, 2 and 3), CG is the category of customers (considering two categories, denoted by *Educ* and *Ret*, standing for educational and retired customers, respectively), A is the age of customers (considering three discretized values, denoted by *Y* (young), *M* (middle) and *O* (old)), C is the city where transactions have been issued (considering three cities, denoted by *NY* (New York), *LA* (Los Angeles) and *SF* (San Francisco)), P is the product of the transactions (considering four products, denoted by c, m, p and r), and Q stands for the quantity of products in the transactions (considering nine such quantities).

Fig. 1 shows the table T in which, for instance, the first tuple means that, at date 1, educational young customers bought 50 products c in New York. Let us now assume that we want to extract all multidimensional sequences that deal with the age of

customers, the products they bought and the corresponding quantities, and that are frequent with respect to the groups of customers and the cities where transactions have been issued. To this end, we consider three sets of dimensions as follows: (i) the dimension D , representing the date, (ii) the three dimensions A , P and Q that we call *analysis dimensions*, (iii) the two dimensions CG and C , that we call *reference dimensions*.

Tuples over analysis dimensions are those that appear in the items that constitute the sequential patterns to be mined. The table is partitioned into blocks according to tuple values over reference dimensions and the support of a given multidimensional sequence is the ratio of the number of blocks supporting the sequence over the total number of blocks. Fig. 2 displays the corresponding blocks in our example.

In this framework, $\{\{(Y, c, 50), (M, p, 2)\}, \{(M, r, 10)\}\}$ is a multidimensional sequence having support $\frac{1}{3}$, since the partition according to the reference dimensions contains 3 blocks, among which one supports the sequence. This is so because $(Y, c, 50)$ and $(M, p, 2)$ both appear at the same date (namely date 1), and $(M, r, 10)$ appears later on (namely at date 2) in the first block shown in Figure 4.

It is important to note that, in our approach, more general patterns, called *jokerized sequences*, can be mined. The reason for this generalization is that considering partially instantiated tuples in sequences implies that more frequent sequences are mined. To see this, considering a support threshold of $\frac{2}{3}$, no sequence of the form $\{\{(Y, c, \mu)\}, \{(M, r, \mu')\}\}$ is frequent. On the other hand, in the first two blocks of Fig. 2, Y associated with c and M associated with r appear one after the other, according to the date of transactions. Thus, we consider that the jokerized sequence, denoted by $\{\{(Y, c, *)\}, \{(M, r, *)\}\}$, is frequent since its support is equal to $\frac{2}{3}$.

D (Date)	CG (Customer-Group)	C (City)	A (Age)	P (Product)	Q (Quantity)
1	Educ	NY	Y	c	50
1	Educ	NY	M	p	2
1	Educ	LA	Y	c	30
1	Ret.	SF	O	c	20
1	Ret.	SF	O	m	2
2	Educ	NY	M	p	3
2	Educ	NY	M	r	10
2	Educ	LA	Y	c	20
3	Educ	LA	M	r	15

Fig. 1. Table T

3 Related Work

In this section, we argue that our approach generalizes previous works on sequential patterns. In particular, the work described in [8] is said to be *intra*-pattern since sequences are mined within the framework of a single description (the so-called *pattern*). In this paper, we propose to generalize this work to *inter*-pattern multidimensional sequences.

3.1 Sequential Patterns

An early example of research in the discovering of patterns from sequences of events can be found in [5]. In this work, the idea is the discovery of rules underlying the generation of a given sequence in order to predict a plausible sequence continuation. This idea is then extended to the discovery of interesting patterns (or *rules*) embedded in a database of sequences of sets of events (items). A more formal approach in solving the problem of mining sequential patterns is the AprioriAll algorithm as presented in [6]. Given a database of sequences, where each sequence is a list of transactions ordered by transaction time, and each transaction is a set of items, the goal is to discover all sequential patterns with a user-specified minimum support, where the support of a pattern is the number of data-sequences that contain the pattern.

In [1], the authors introduce the problem of mining sequential patterns over large databases of customer transactions where each transaction consists of customer-id, transaction time, and the items bought in the transaction. Formally, given a set of sequences, where each sequence consists of a list of elements and each element consists of a set of items, and given a user-specified min support threshold, sequential pattern mining is to find all of the frequent subsequences, i.e., the subsequences whose occurrence frequency in the set of sequences is no less than min support. Sequential pattern mining discovers frequent patterns ordered by time. An example of this type of pattern is *A customer who bought a new television 3 months ago, is likely to buy a DVD player now*. Subsequently, many studies have introduced various methods in mining sequential patterns (mainly in time-related data) but almost all proposed methods are Apriori-like, i.e., based on the Apriori property which states the fact that any super-pattern of a nonfrequent pattern cannot be frequent. An example using this approach is the GSP algorithm [9].

3.2 Multidimensional Sequential Patterns

As far as we know, three propositions have been studied in order to deal with several dimensions when building sequential patterns. Next, we briefly recall these propositions.

Pinto et al. [8]. This work is the first one dealing with several dimensions in the framework of sequential patterns. For instance, purchases are not only described by considering the customer ID and the products, but also by considering the age, the type of the customer (Cust-Grp) and the city where (s)he lives, as shown in Fig. 1.

Multidimensional sequential patterns are defined over the schema A_1, \dots, A_m, S where A_1, \dots, A_m are the dimensions describing the data and S is the sequence of items purchased by the customers ordered over time. A multidimensional sequential pattern is defined as $(id_1, (a_1, \dots, a_m), s)$ where $a_i \in A_i \cup \{*\}$. $id_1, (a_1, \dots, a_m)$ is said to be a multidimensional pattern. For instance, the authors consider the sequence $((*, NY, *), \langle bf \rangle)$ meaning that customers from NY have all bought a product b and then a product f . Sequential patterns are mined from such multidimensional databases either (i) by mining all frequent sequential patterns over the product dimension and then regrouping them into multidimensional patterns, (ii) or by mining all frequent multidimensional patterns and then mining frequent product sequences over these patterns. Note that the sequences found by this approach do not contain several dimensions since the dimension time only

concerns products. Dimension product is the only dimension that can be combined over time, meaning that it is not possible to have a rule indicating that when b is bought in *Boston* then c is bought in *NY*. Therefore, our approach can be seen as a generalization of the work in [8].

Yu et Chen. [11]. In this work, the authors consider sequential pattern mining in the framework of Web Usage Mining. Even if three dimensions (pages, sessions, days) are considered, these dimensions are very particular since they belong to a single hierarchized dimension. Thus, the sequences mined in this work describe correlations between objects over time by considering only one dimension, which corresponds to the web pages.

de Amo et al. [4]. This approach is based on first order temporal logic. This proposition is close to our approach, but more restricted since (i) groups used to compute the support are predefined whereas we consider the fact that the user should be able to define them (see reference dimensions below), and (ii) several attributes cannot appear in the sequences. The authors claim that they aim at considering several dimensions but they have only shown one dimension for the sake of simplicity. However, the paper does not provide hints for a complete solution with *real* multidimensional patterns, as we do in our approach.

4 M²SP: Mining Multidimensional Sequential Patterns

4.1 Dimension Partition

For each table defined on the set of dimensions D , we consider a partition of D into four sets: D_t for the temporal dimension, D_A for the *analysis* dimensions, D_R for the *reference* dimensions, and D_F for the *ignored* dimensions.

Each tuple $c = (d_1, \dots, d_n)$ can thus be written as $c = (f, r, a, t)$ where f, r, a and t are the restrictions of c on D_F, D_R, D_A and D_t , respectively.

Given a table T , the set of all tuples in T having the same restriction r over D_R is said to be a *block*. Each such block B is denoted by the tuple r that defines it, and we denote by B_{T, D_R} the set of all blocks that can be built up from table T .

In our running example, we consider $F = \emptyset$, $D_R = \{CG, C\}$, $D_A = \{A, P, Q\}$ and $D_t = \{D\}$. Fig. 2 shows the three blocks built up from table T .

D	CG	C	A	P	Q	D	CG	C	A	P	Q	D	CG	C	A	P	Q
1	Educ	NY	Y	c	50	1	Educ	LA	Y	c	30	1	Ret.	SF	O	c	20
1	Educ	NY	M	p	2	2	Educ	LA	Y	c	20	1	Ret.	SF	O	m	2
2	Educ	NY	M	p	3	3	Educ	LA	M	r	15						
2	Educ	NY	M	r	10												

a. Block (*Educ, NY*)

b. Block (*Educ, LA*)

c. Block (*Ret., SF*)

Fig. 2. Blocks defined on T over dimensions CG and C

When mining multidimensional sequential patterns, the set D_R identifies the blocks of the database to be considered when computing supports. The support of a sequence is the proportion of blocks embedding it. Note that, in the case of usual sequential patterns and of sequential patterns as in [8] and [4], this set is reduced to one dimension (*cid* in [8] or *IdG* in [4]).

The set D_A describes the analysis dimensions, meaning that values over these dimensions appear in the multidimensional sequential patterns. Note that usual sequential patterns only consider one analysis dimension corresponding to the products purchased or the web pages visited. The set F describes the ignored dimensions, i.e. those that are used neither to define the date, nor the blocks, nor the patterns to be mined.

4.2 Multidimensional Item, Itemset and Sequential Pattern

Definition 1 (Multidimensional Item). Let $D_A = \{D_{i_1}, \dots, D_{i_m}\}$ be a subset of D . A multidimensional item on D_A is a tuple $e = (d_{i_1}, \dots, d_{i_m})$ such that, for every k in $[1, m]$, d_{i_k} is in $Dom(D_{i_k})$.

Definition 2 (Multidimensional Itemset). A multidimensional itemset on D_A is a non empty set of items $i = \{e_1, \dots, e_p\}$ where for every j in $[1, p]$, e_j is a multidimensional item on D_A and for all j, k in $[1, p]$, $e_j \neq e_k$.

Definition 3 (Multidimensional Sequence). A multidimensional sequence on D_A is an ordered non empty list of itemsets $\zeta = \langle i_1, \dots, i_l \rangle$ where for every j in $[1, l]$, i_j is a multidimensional itemset on D_A .

In our running example, $(Y, c, 50)$, $(M, p, 2)$, $(M, r, 10)$ are three multidimensional items on $D_A = \{A, P, Q\}$. Thus, $\langle \{(Y, c, 50), (M, p, 2)\}, \{(M, r, 10)\} \rangle$ is a multidimensional sequence on D_A .

Definition 4 (Inclusion of sequence). A multidimensional sequence $\zeta = \langle a_1, \dots, a_l \rangle$ is said to be a subsequence of a sequence $\zeta' = \langle b_1, \dots, b_{l'} \rangle$ if there exist $1 \leq j_1 \leq j_2 \leq \dots \leq j_l \leq l'$ such that $a_1 \subseteq b_{j_1}, a_2 \subseteq b_{j_2}, \dots, a_l \subseteq b_{j_l}$.

With $\zeta = \langle \{(Y, c, 50)\}, \{(M, r, 10)\} \rangle$ and $\zeta' = \langle \{(Y, c, 50), (M, p, 2)\}, \{(M, r, 10)\} \rangle$, ζ is a subsequence of ζ' .

4.3 Support

Computing the support of a sequence amounts to count the number of blocks that *support* the sequence. Intuitively, a block supports a sequence ζ if (i) for each itemset i in ζ there exists a date in $Dom(D_t)$ such that all items in i appear at this date, and (ii) all itemsets in ζ are successively retrieved at different and increasing dates.

Definition 5. A table T supports a sequence $\langle i_1, \dots, i_l \rangle$ if for every $j = 1, \dots, l$, there exists d_j in $Dom(D_t)$ such that for every item e in i_j , there exists $t = (f, r, e, d_j)$ in T with $d_1 < d_2 < \dots < d_l$.

In our running example, the block $(Educ, NY)$ from Fig. 2.a supports $\zeta = \langle \{(Y, c, 50), (M, p, 2)\}, \{(M, r, 10)\} \rangle$ since $\{(Y, c, 50), (M, p, 2)\}$ appears at $date = 1$ and $\{(M, r, 10)\}$ appears at $date = 2$.

The support of a sequence in a table T is the proportion of blocks of T that support it.

Definition 6 (Sequence Support). Let D_R be the reference dimensions and T a table partitioned into the set of blocks B_{T, D_R} . The support of a sequence ζ is defined by:

$$support(\zeta) = \frac{|\{B \in B_{T, D_R} \mid B \text{ supports } \zeta\}|}{|B_{T, D_R}|}$$

Definition 7 (Frequent Sequence). Let $minsup \in [0, 1]$ be the minimum user-defined support value. A sequence ζ is said to be frequent if $support(\zeta) \geq minsup$. An item e is said to be frequent if so is the sequence $\langle \{e\} \rangle$.

In our running example, let us consider $D_R = \{CG, C\}$, $D_A = \{A, P, Q\}$, $minsup = \frac{1}{5}$, $\zeta = \langle \{(Y, c, 50), (M, p, 2)\}, \{(M, r, 10)\} \rangle$. The three blocks of the partition of T from Fig. 2 must be scanned to compute $support(\zeta)$.

1. Block (Educ, NY) (Fig. 2.a). In this block, we have $(Y, c, 50)$ and $(M, p, 2)$ at date 1, and $(M, r, 10)$ at date 2. Thus this block supports ζ .

2. Block (Educ, LA) (Fig. 2.b). This block does not support ζ since it does not contain $(M, p, 2)$.

3. Block (Ret., SF) (Fig. 2.c). This block does not support ζ since it contains only one date.

Thus, we have $support(\zeta) = \frac{1}{3} \geq minsup$.

5 Jokerized Sequential Patterns

Considering the definitions above, an item can only be retrieved if there exists a frequent tuple of values from domains of D_A containing it. For instance, it can happen that neither (Y, r) nor (M, r) nor (O, r) is frequent whereas the value r is frequent. In this case, we consider $(*, r)$ which is said to be *jokerized*.

Definition 8 (Jokerized Item). Let $e = (d_1, \dots, d_m)$ a multidimensional item. We denote by $e_{[d_i/\delta]}$ the replacement in e of d_i by δ . e is said to be a *jokerized multidimensional item* if: (i) $\forall i \in [1, m], d_i \in Dom(D_i) \cup \{*\}$, and (ii) $\exists i \in [1, m]$ such that $d_i \neq *$, and (iii) $\forall d_i = *, \nexists \delta \in Dom(D_i)$ such that $e_{[d_i/\delta]}$ is frequent.

A *jokerized* item contains at least one specified analysis dimension. It contains a $*$ only if no specific value from the domain can be set. A *jokerized* sequence is a sequence containing at least one *jokerized* item. A block is said to *support* a sequence if a set of tuples containing the itemsets satisfying the temporal constraints can be found.

Definition 9 (Support of a Jokerized Sequence). A table T supports a *jokerized* sequence $\zeta = \langle i_1, \dots, i_l \rangle$ if: $\forall j \in [1, l], \exists \delta_j \in Dom(D_{t_j}), \forall e = (d_{i_1}, \dots, d_{i_m}) \in i_j, \exists t = (f, r, (x_{i_1}, \dots, x_{i_m}), \delta_j) \in T$ with $d_{i_k} = x_{i_k}$ or $d_{i_k} = *$ and $\delta_{i_1} < \delta_{i_2} < \dots < \delta_{i_l}$.

The support of ζ is defined by: $support(\zeta) = \frac{|\{B \in B_{T, D_R} \text{ s.t. } B \text{ supports } \zeta\}|}{|B_{T, D_R}|}$

6 Algorithms

6.1 Mining Frequent Items

The computation of all frequent sequences is based on the computation of all frequent multidimensional items. When considering no joker value, a single scan of the database is enough to compute them.

On the other hand, when considering jokerized items, a levelwise algorithm is used in order to build the frequent multidimensional items having as few joker values as possible. To this end, we consider a lattice which lower bound is the multidimensional item $(*, \dots, *)$. This lattice is partially built from $(*, \dots, *)$ up to the frequent items containing as few $*$ as possible. At level i , i values are specified, and items at this level are combined to build a set of candidates at level $i+1$. Two frequent items are combined to build a candidate if they are \bowtie -compatible.

Definition 10 (\bowtie -compatibility). Let $e_1 = (d_1, \dots, d_n)$ and $e_2 = (d'_1, \dots, d'_n)$ be two distinct multidimensional items where d_i and $d'_i \in \text{dom}(D_i) \cup \{*\}$. e_1 and e_2 are said to be \bowtie -compatible if there exists $\Delta = \{D_{i_1}, \dots, D_{i_{n-2}}\} \subset \{D_1, \dots, D_n\}$ such that for every $j \in [1, n-2]$, $d_{i_j} = d'_{i_j} \neq *$ with $d_{i_{n-1}} = *$ and $d'_{i_{n-1}} \neq *$ and $d_{i_n} \neq *$ and $d'_{i_n} = *$.

Definition 11 (Join). Let $e_1 = (d_1, \dots, d_n)$ and $e_2 = (d'_1, \dots, d'_n)$ be two \bowtie -compatible multidimensional items. We define $e_1 \bowtie e_2 = (v_1, \dots, v_n)$ where $v_i = d_i$ if $d_i = d'_i$, $v_i = d_i$ if $d'_i = *$ and $v_i = d'_i$ if $d_i = *$.

Let E and E' be two sets of multidimensional items of size n , we define

$$E \bowtie E' = \{e \bowtie e' \mid (e, e') \in E \times E' \wedge e \text{ and } e' \text{ are } \bowtie\text{-compatible}\}$$

In our running example, $(NY, Y, *)$ and $(*, Y, r)$ are \bowtie -compatible. We have $(NY, Y, *) \bowtie (*, Y, r) = (NY, Y, r)$. On the contrary, $(NY, M, *)$ and $(NY, Y, *)$ are not \bowtie -compatible. Note that this method is close to the one used for iceberg cubes in [2,3].

Let F_1^i denote the set of 1-frequent items having i dimensions which are specified (different from $*$). F_1^1 is obtained by counting each value over each analysis dimension, i.e., $F_1^1 = \{f \in \text{Cand}_1^1, \text{support}(f) \geq \text{minsup}\}$. Candidate items of size i are obtained by joining the set of frequent items of size $i-1$ with itself: $\text{Cand}_1^i = F_1^{i-1} \bowtie F_1^{i-1}$.

Function supportcount

Data : ζ, T, D_R , counting //counting indicates if joker values are considered or not

Result : support of ζ

Integer support $\leftarrow 0$; Boolean seqSupported;

$B_{T, D_R} \leftarrow \{\text{blocks of } T \text{ identified over } D_R\}$;

foreach $B \in B_{T, D_R}$ **do**

 seqSupported $\leftarrow \text{supportTable}(\zeta, B, \text{counting})$;

if seqSupported **then** support $\leftarrow \text{support} + 1$;

return $\left(\frac{\text{support}}{|B_{T, D_R}|} \right)$

Algorithm 1: Support of a sequence (supportcount)

Function supportTable**Data** : $\zeta, T, \text{counting}$ **Result** : Boolean $ItemSetFound \leftarrow false$; $seq \leftarrow \zeta$; $itset \leftarrow seq.first()$; $it \leftarrow itset.first()$ **if** $\zeta = \emptyset$ **then** **return** (true) // End of Recursivity**while** $t \leftarrow T.next \neq \emptyset$ **do** **if** $supports(t, it, \text{counting})$ **then** **if** $(NextItem \leftarrow itset.second()) = \emptyset$ **then** $ItemSetFound \leftarrow true$

// Look for all the items from the itemset

else

// Anchoring on the item (date)

 $T' \leftarrow \sigma_{date=t.date}(T)$ **while** $t' \leftarrow T'.next() \neq \emptyset \wedge ItemSetFound = false$ **do** **if** $supports(t', NextItem, \text{counting})$ **then** $NextItem \leftarrow itset.next()$ **if** $NextItem = \emptyset$ **then** $ItemSetFound \leftarrow true$ **if** $ItemSetFound = true$ **then** // Anchoring on the current itemset succeeded; test the other itemsets in seq **return** ($supportTable(seq.tail(), \sigma_{date>t.date}(T), \text{counting})$) **else**

// Anchoring failure: try anchoring with the next dates

 $itset \leftarrow seq.first()$ $T \leftarrow \sigma_{date>t.date}(T)$ // Skip to next dates**return**(false) // Not found**Algorithm 2: supportTable** (Checks if a sequence ζ is supported by a table T)

6.2 Mining Jokerized Multidimensional Sequences

The frequent items give all frequent sequences containing one itemset consisting of a single item. Then, the candidate sequences of size k ($k \geq 2$) are generated and validated against the table T . This computation is based on usual algorithms such as PSP [7] that are adapted for the treatment of joker values.

The computation of the support of a sequence ζ according to the reference dimensions D_R is given by Algorithm 1. This algorithm checks whether each block of the partition supports the sequence by calling the function supportTable (Algorithm 2). *supportTable* attempts to find a tuple from the block that matches the first item of the first itemset of the sequence in order to *anchor* the sequence. This operation is repeated recursively until all itemsets from the sequence are found (return true) or until there is no way to go on further (return false). Several possible anchors may have to be tested.

7 Experiments

In this section, we report experiments performed on synthetic data. These experiments aim at showing the interest and scalability of our approach, especially in the jokerized approach. As many databases from the real world include quantitative information, we

have distinguished a quantitative dimension. In order to highlight the particular role of this quantitative dimension, we consider four ways of computing frequent sequential patterns: (i) no joker (M^2SP), (ii) jokers on all dimensions but the quantitative one ($M^2SP-alpha$), (iii) jokers only on the quantitative dimension (M^2SP-mu), (iv) jokers on all dimensions ($M^2SP-alpha-mu$). Note that case (iv) corresponds to the jokerized approach presented in Section 5. Our experiments can thus be seen as being conducted in the context of a fact table of a multidimensional database, where the quantitative dimension is the *measure*. In Figures 5-12, minsup is the minimum support taken into account, nb_dim is the number of analysis dimensions being considered, DB_size is the number of tuples, and avg_card is the average number of values in the domains of the analysis dimensions.

Fig. 3 and 4 compare the behavior of the four approaches described above when the support changes. $M^2SP-alpha$ and $M^2SP-alpha-mu$ have a similar behavior, the difference being due to the verification of quantities in the case of $M^2SP-alpha$. Note that these experiments are not led with the same minimum support values, since no frequent items are found for M^2SP and M^2SP-mu if the support is too high. Fig. 5 shows the scalability of our approach since runtime grows almost linearly when the database size increases (from 1,000 tuples up to 26,000 tuples).

Fig. 6 shows how runtime behaves when the average cardinality of the domains of analysis dimensions changes. When this average is very low, numerous frequent items are mined among few candidates. On the contrary, when this average is high, numerous candidates have to be considered from which few frequent items are mined. Between these two extrema, the runtime decreases. Fig. 7 and 8 show the behavior of our approach when the number of analysis dimensions changes. The number of frequent items increases as the number of analysis dimensions grows, leading to an increase of the number of frequent sequences. Fig. 9 and 10 show the differential between the number of frequent sequences mined by our approach compared to the number of frequent sequences mined by the approach described in [8], highlighting the interest of our proposition.

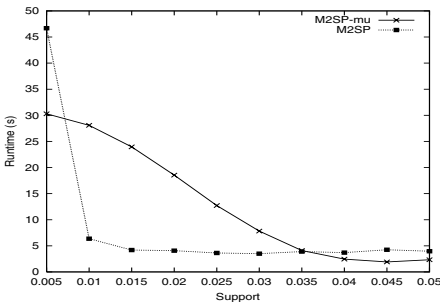


Fig. 3. Runtime over Support (DB_size=12000, nb_dim=5, avg_card=20)

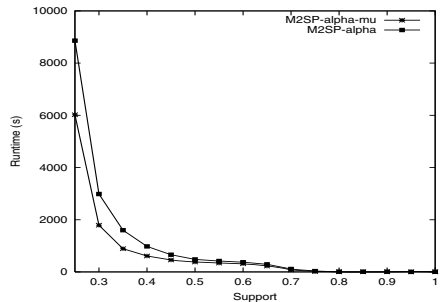


Fig. 4. Runtime over Support (DB_size=12000, nb_dim=5, avg_card=20)

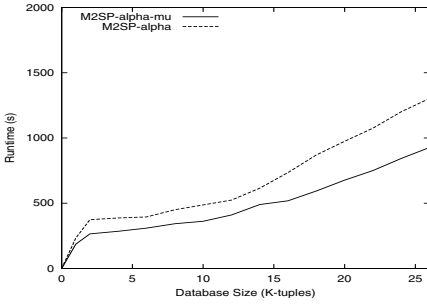


Fig. 5. Runtime over database size (minsup=0.5, nb_dim=15, avg_card = 20)

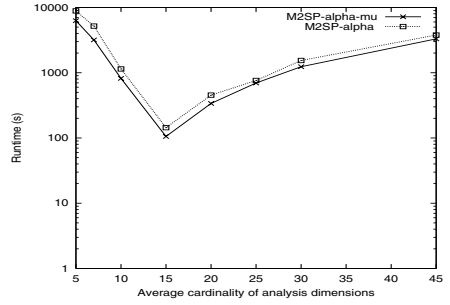


Fig. 6. Runtime over Average Cardinality of Analysis Dimensions (minsup=0.8, DB_size=12000, nb_dim=15)

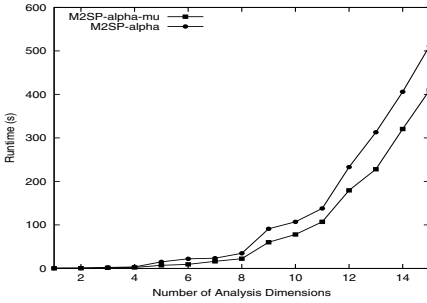


Fig. 7. Runtime over Number of Analysis Dimensions (minsup=0.5, DB_size=12000, nb_dim=15, avg_card=20)

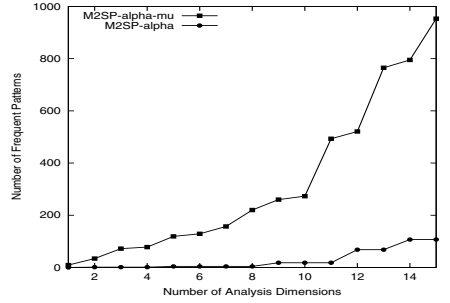


Fig. 8. Number of Frequent patterns over number of analysis dimensions (minsup=0.5, DB_size=12000, nb_dim=15, avg_card=20)

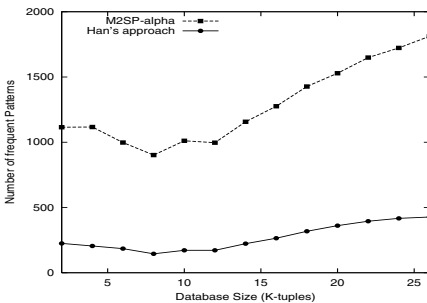


Fig. 9. Number of Frequent Sequences over Database Size (minsup=0.5, nb_dim=15, avg_card=20)

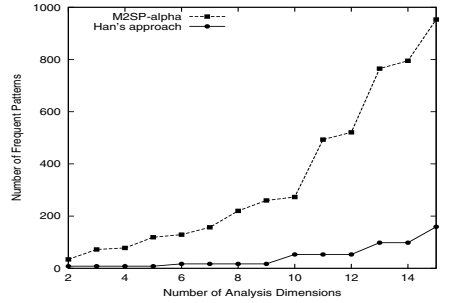


Fig. 10. Number of Frequent Sequences over Number of Analysis Dimensions (minsup=0.5, DB_size=12000, avg_card=20)

8 Conclusion

In this paper, we have proposed a novel definition for multidimensional sequential patterns. Contrary to the propositions [4,8,11], several analysis dimensions can be found in the sequence, which allows for the discovery of rules as *A customer who bought a surfboard together with a bag in NY later bought a wetsuit in LA*. We have also defined *jokerized sequential patterns* by introducing the joker value * on analysis dimensions. Algorithms have been evaluated against synthetic data, showing the scalability of our approach.

This work can be extended following several directions. For example, we can take into account approximate values on quantitative dimensions. In this case, we allow the consideration of values that are not fully jokerized while remaining frequent. This proposition is important when considering data from the real world where the high number of quantitative values prevents each of them to be frequent. Rules to be built will then be like *The customer who bought a DVD player on the web is likely to buy almost 3 DVDs in a supermarket later*. Hierarchies can also be considered in order to mine multidimensional sequential patterns at different levels of granularity in the framework of multidimensional databases.

References

1. R. Agrawal and R. Srikant. Mining sequential patterns. In *Proc. 1995 Int. Conf. Data Engineering (ICDE'95)*, pages 3–14, 1995.
2. K. Beyer and R. Ramakrishnan. Bottom-up computation of sparse and iceberg cube. In *Proc. of ACM SIGMOD Int. Conf. on Management of Data*, pages 359–370, 1999.
3. A. Casali, R. Cicchetti, and L. Lakhal. Cube lattices: A framework for multidimensional data mining. In *Proc. of 3rd SIAM Int. Conf. on Data Mining*, 2003.
4. S. de Amo, D. A. Furtado, A. Giacometti, and D. Laurent. An apriori-based approach for first-order temporal pattern mining. In *Simpósio Brasileiro de Bancos de Dados*, 2004.
5. T.G. Dietterich and R.S. Michalski. Discovering patterns in sequences of events. *Artificial Intelligence*, 25(2):187–232, 1985.
6. H. Mannila, H. Toivonen, and A.I. Verkamo. Discovering frequent episodes in sequences. In *Proc. of Int. Conf. on Knowledge Discovery and Data Mining*, pages 210–215, 1995.
7. F. Massegli, F. Cathala, and P. Poncelet. The PSP Approach for Mining Sequential Patterns. In *Proc. of PKDD*, volume 1510 of *LNCS*, pages 176–184, 1998.
8. H. Pinto, J. Han, J. Pei, K. Wang, Q. Chen, and U. Dayal. Multi-dimensional sequential pattern mining. In *ACM CIKM*, pages 81–88, 2001.
9. R. Srikant and R. Agrawal. Mining Sequential Patterns: Generalizations and Performance Improvements. In *Proc. of EDBT*, pages 3–17, 1996.
10. J.D. Ullman. *Principles of Database and Knowledge-Base Systems*, volume I. Computer Science Press, 1988.
11. C.-C. Yu and Y.-L. Chen. Mining sequential patterns from multidimensional sequence data. *IEEE Transactions on Knowledge and Data Engineering*, 17(1):136–140, 2005.

A Systematic Comparison of Feature-Rich Probabilistic Classifiers for NER Tasks

Benjamin Rosenfeld, Moshe Fresko, and Ronen Feldman

Bar-Ilan University, Computer Science Department, Data-Mining Lab.,
Ramat-Gan, Israel
{freskom1, feldman}@cs.biu.ac.il
<http://www.cs.biu.ac.il/~freskom1/DataMiningLab>

Abstract. In the CoNLL 2003 NER shared task, more than two thirds of the submitted systems used the feature-rich representation of the task. Most of them used maximum entropy to combine the features together. Others used linear classifiers, such as SVM and RRM. Among all systems presented there, one of the MEMM-based classifiers took the second place, losing only to a committee of four different classifiers, one of which was ME-based and another RRM-based. The lone RRM was fourth, and CRF came in the middle of the pack. In this paper we shall demonstrate, by running the three algorithms upon the same tasks under exactly the same conditions that this ranking is due to feature selection and other causes and not due to the inherent qualities of the algorithms, which should be ranked otherwise.

1 Introduction

Recently, feature-rich probabilistic conditional classifiers became state-of-the-art in sequence labeling tasks, such as NP chunking, PoS tagging, and Named Entity Recognition. Such classifiers build a probabilistic model of the task, which defines a conditional probability on the space of all possible labelings of a given sequence. In this, such classifiers differ from the binary classifiers, such as decision trees and rule-based systems, which directly produce classification decisions, and from the generative probabilistic classifiers, such as HMM-based Nymble [2] and SCFG-based TEG [8], which model the joint probability of sequences and their labelings. Modeling the conditional probability allows the classifiers to have all the benefits of probabilistic systems while having the ability to use any property of tokens and their contexts, if the property can be represented in the form of binary features. Since almost all local properties can be represented in such a way, this ability is very powerful.

There are several different feature-rich probabilistic classifiers developed by different researchers, and in order to compare them, one usually takes a known publicly available dataset, such as MUC-7 [23] or CoNLL shared task [12], and compares the performance of the algorithms on the dataset. However, performance of a feature-rich classifier strongly depends upon the feature sets it uses. Since systems developed by different researches are bound to use different feature sets, the differences in performance of complete systems can not reliably teach us about the qualities of the underlying algorithms.

In this work we compare the performances of three common models (all present in the CoNLL 2003 shared task) – MEMM [15], CRF [16], and RRM (regularized Winnow) [14] – within the same platform, using exactly the same set of features. We also test the effects of different training sizes, different choice of parameters, and different feature sets upon the algorithms' performance.

Our experiments indicate that CRF outperforms MEMM for all datasets and feature sets, which is not surprising, since CRF is a better model of sequence labeling. Surprisingly, though, the RRM performs at the same level or even better than CRF, despite being local model like MEMM, and being significantly simpler to build than both CRF and MEMM.

The following section of the paper we outline the three algorithms. We then present our experiments and their results.

2 Classifiers

The general sequence labeling problem can be described as follows. Given a small set Y of labels, and a sequence $\mathbf{x} = x_1x_2\dots x_{l(\mathbf{x})}$, the task is to find a labeling $\mathbf{y} = y_1y_2\dots y_{l(\mathbf{x})}$, where each label $y_i \in Y$. In the framework of feature-rich classifiers, the elements x_i of the sequence should not be thought of as simple tokens, but rather as sequence positions, or *contexts*. The contexts are characterized by a set of externally supplied binary *features*. Thus, each context x_i can be represented as a vector $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ik})$, where $x_{ij} = 1$ if the j -th feature is present in the i -th context, and $x_{ij} = 0$ otherwise.

The feature-rich sequence classifiers have no knowledge of the nature of features and labels. Instead, in order to make predictions, the classifiers are supplied with a training set $T = \{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})\}_{t=1..n}$ of sequences with their intended labelings. The classifiers use the training set to build the model of the task, which is subsequently used to label unseen sequences.

We shall describe the particular algorithms only briefly, referring to the original works to supply the details.

2.1 MEMM

A Maximum Entropy Markov Model classifier [4] builds a probabilistic conditional model of sequence labeling. Labeling each position in each sequence is considered to be a separate classification decision, possibly influenced by a small constant number of previous decisions in the same sequence. In our experiments we use a Markov model of order one, in which only the most recent previous decision is taken into account.

Maximal Entropy models are formulated in terms of *feature functions* $f(\mathbf{x}_i, y_i, y_{i-1}) \rightarrow \{0, 1\}$, which link together the context features and the target labels. In our formulation, we have a feature function f_{jy} for each context feature j and each label y , and a feature function $f_{jyy'}$ for each context feature and each pair of labels. The functions are defined as follows:

$f_{jy}(\mathbf{x}_i, y_i, y_{i-1}) = \mathbf{x}_{ij}I_y(y_i)$ and $f_{jyy'}(\mathbf{x}_i, y_i, y_{i-1}) = \mathbf{x}_{ij}I_y(y_i)I_{y'}(y_{i-1})$, where $I_a(b)$ is one if $a = b$ and zero otherwise. The vector of all feature functions is denoted $\mathbf{f}(\mathbf{x}_i, y_i, y_{i-1})$.

A trained MEMM model has a real weight λ_f for each feature function f . Together, the weights form the parameter vector $\boldsymbol{\lambda}$. The model has the form

$$(1) \quad P_{\boldsymbol{\lambda}}(y_i \mid \mathbf{x}_i, y_{i-1}) = \frac{1}{Z(\mathbf{x}_i, y_{i-1})} \exp(\boldsymbol{\lambda} \cdot \mathbf{f}(\mathbf{x}_i, y_i, y_{i-1})),$$

where $Z(\mathbf{x}_i, y_{i-1}) = \sum_y \exp(\boldsymbol{\lambda} \cdot \mathbf{f}(\mathbf{x}_i, y, y_{i-1}))$ is the factor making the probabilities for different labels sum to one.

Given a model (1), it can be used for inferring the labeling $\mathbf{y} = y_1 y_2 \dots y_{l(x)}$ of an unseen sequence $\mathbf{x} = \mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_{l(x)}$ by calculating the most probable overall sequence of labels:

$$(2) \quad \mathbf{y}(\mathbf{x}) := \arg \max_{y_1 y_2 \dots y_{l(x)}} \sum_{i=1}^{l(x)} \log P_{\boldsymbol{\lambda}}(y_i \mid \mathbf{x}_i, y_{i-1}).$$

This most probable sequence can be efficiently calculated using a variant of the Viterbi algorithm.

The model parameters are trained in such a way as to maximize the model’s entropy while making the expected value of each feature function agree with the observed relative frequency of the feature function in the training data. Those conditions can be shown to be uniquely satisfied by the model which maximizes the log-likelihood of the training data among all models of the form (1). In order to avoid overfitting, the likelihood can be penalized with a prior $\Pr(\boldsymbol{\lambda})$. Then, the log-likelihood is

$$\begin{aligned} L_T(\boldsymbol{\lambda}) &= \sum_t \sum_{i=1}^{l(x^{(t)})} \log P_{\boldsymbol{\lambda}}(y_i^{(t)} \mid \mathbf{x}_i^{(t)}, y_{i-1}^{(t)}) - \Pr(\boldsymbol{\lambda}) = \\ &= \sum_t \sum_{i=1}^{l(x^{(t)})} (\boldsymbol{\lambda} \cdot \mathbf{f}(\mathbf{x}_i^{(t)}, y_i^{(t)}, y_{i-1}^{(t)}) - \log Z(\mathbf{x}_i^{(t)})) - \Pr(\boldsymbol{\lambda}) \end{aligned}$$

and its gradient is

$$\nabla L_T(\boldsymbol{\lambda}) = \sum_t \sum_{i=1}^{l(x^{(t)})} (\mathbf{f}(\mathbf{x}_i^{(t)}, y_i^{(t)}, y_{i-1}^{(t)}) - E_{P_{\boldsymbol{\lambda}}}(\mathbf{f}(\mathbf{x}_i^{(t)}, Y, y_{i-1}^{(t)}))) - \nabla \Pr(\boldsymbol{\lambda}),$$

where

$$E_{P_{\boldsymbol{\lambda}}}(\mathbf{f}(\mathbf{x}_i^{(t)}, Y, y_{i-1}^{(t)})) = \sum_{y \in Y} P_{\boldsymbol{\lambda}}(y \mid \mathbf{x}_i^{(t)}, y_{i-1}^{(t)}) \mathbf{f}(\mathbf{x}_i^{(t)}, y, y_{i-1}^{(t)})$$

is the expectation of the feature vector under the model (1).

With a reasonably chosen prior, the function $L_T(\boldsymbol{\lambda})$ is strictly concave, and so can be maximized by any convex optimization algorithm. We use L-BFGS for this purpose.

2.2 CRF

A Conditional Random Fields (CRF) [7] classifier also builds a probabilistic model of sequence labeling. CRF uses the maximal entropy principle to model the labeling of a sequence as a whole, in contrast to MEMM, which builds a model of separate labeling decisions at different sequence positions.

The model is built upon exactly the same vector $\mathbf{f}(\mathbf{x}_i, y_i, y_{i-1})$ of feature functions as MEMM. The feature functions are summed along a sequence to produce a *sequence feature functions vector*

$$(3) \quad \mathbf{F}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{l(\mathbf{x})} \mathbf{f}(\mathbf{x}_i, y_i, y_{i-1}),$$

which is then used for constructing the maximal entropy model

$$P_{\lambda}(y | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\lambda \cdot \mathbf{F}(\mathbf{x}, y)).$$

A trained model can be used for inferring the most probable labeling of an unseen sequence. The decomposition (3) allows to use the Viterbi algorithm almost identically to the MEMM case, except that in (2), instead of $\log P_{\lambda}(y_i | \mathbf{x}_i, y_{i-1}) = \lambda \cdot \mathbf{f}(\mathbf{x}_i, y_i, y_{i-1}) - \log Z(\mathbf{x}_i, y_{i-1})$, simple $\lambda \cdot \mathbf{f}(\mathbf{x}_i, y_i, y_{i-1})$ is used. Since $Z(\mathbf{x})$ does not depend on labeling, it need not be calculated at all during inference.

To train the CRF model, we need to maximize the model entropy while satisfying the expectation constrains, expressed this time in terms of the sequence feature functions. As before, this is equivalent to maximizing the log-likelihood of the training data, which can also be penalized with a prior to avoid overfitting:

$$L_T(\lambda) = \sum_t \log P_{\lambda}(y^{(t)} | \mathbf{x}^{(t)}) - \frac{\|\lambda\|^2}{2\sigma^2} = \sum_t (\lambda \cdot \mathbf{F}(\mathbf{x}^{(t)}, y^{(t)}) - \log Z(\mathbf{x}^{(t)})) - \Pr(\lambda).$$

The gradient is

$$\nabla L_T(\lambda) = \sum_t (\mathbf{F}(\mathbf{x}^{(t)}, y^{(t)}) - E_{P_{\lambda}}(\mathbf{F}(\mathbf{x}^{(t)}, \mathbf{Y}^{(t)}))) - \nabla \Pr(\lambda),$$

where $\mathbf{Y}^{(t)}$ is the set of label sequences of length $l(\mathbf{x}^{(t)})$, and $E_{P_{\lambda}}(\mathbf{F}(\mathbf{x}^{(t)}, \mathbf{Y})) = \sum_{y \in \mathbf{Y}} P_{\lambda}(y | \mathbf{x}^{(t)}) \mathbf{F}(\mathbf{x}^{(t)}, y)$

is the expectation of the sequence feature functions vector under the model (3).

In order to maximize $L_T(\lambda)$, we need a way to calculate $\log Z(\mathbf{x})$ and $E_{P_{\lambda}}(\mathbf{F}(\mathbf{x}, \mathbf{Y}))$ for the given sequence \mathbf{x} . It is possible to do this efficiently, using a variant of the Forward-Backward algorithm. Details can be found in [7] and [19].

2.3 RRM

The Robust Risk Minimization classifier [14] results from regularization of the Winnow algorithm [21]. Winnow is a multiplicative-update online algorithm used for estimating the weights of a binary linear classifier, which has the following general form:

$$y = \text{sign}(\mathbf{w}^T \mathbf{x}),$$

where \mathbf{x} is the input vector, \mathbf{w} is the weight vector, and $y \in \{+1, -1\}$ is the classification decision.

It was shown in [20], that using a risk function of a special form, the regularized Winnow can produce such weights \mathbf{w} that

$$P(y = +1 | \mathbf{x}) \approx (Tr_{[-1,1]}(\mathbf{w}^T \mathbf{x}) + 1) / 2,$$

where $Tr_{[a,b]}(s) = \min(b, \max(a, s))$ is a truncation of s onto $[a, b]$.

Although the derivation is elaborate, the resulting algorithm is very simple. It consists of iteratively going over the training set $T = \{(\mathbf{x}^{(t)}, y^{(t)})\}_{t=1..n}$ (here, $y^{(t)} = \pm 1$), and incrementally updating

$$(4) \quad \begin{aligned} \alpha_t &:= Tr_{[0,2c]} \left(\alpha_t + \eta \left(1 - \frac{\alpha_t}{c} - \mathbf{w}^T \mathbf{x}^{(t)} y^{(t)} \right) \right) \\ w_j &:= \mu_j \exp \left(\sum_t \alpha_t x_j^{(t)} y^{(t)} \right) \end{aligned}$$

The α_t are the *dual* weights, initialized to zero and kept between the iterations. c is the *regularization* parameter, η is the *learning rate*, and μ_j is the *prior*.

The $y^{(t)}$ in (4) are binary decisions. In order to use the RRM for sequence labeling task with more than two labels, we can build a separate classifier for each label and then combine them together within a single Viterbi search.

3 Experimental Setup

The goal of this work is to compare the three sequence labeling algorithms in several different dimensions: absolute performance, dependence upon the corpus, dependence upon the training set size and the feature set, and dependence upon the hyper-parameters.

3.1 Datasets

For our experiments we used four datasets: CoNLL-E, the English CoNLL 2003 shared task dataset, CoNLL-D, the German CoNLL 2003 shared task dataset, the MUC-7 dataset [23], and the proprietary CLF dataset [8]. For the experiments with smaller training sizes, we cut training corpora into chunks of 10K, 20K, 40K, 80K, and 160K tokens. The corresponding datasets are denoted $\langle \text{Corpus} \rangle_{\langle \text{Size} \rangle}$, e.g. “CoNLL-E_10K”.

3.2 Feature Sets

There are many properties of tokens and their contexts that could be used in a NER system. We experiment with the following properties, ordered according to the difficulty of obtaining them:

- A. The exact character strings of tokens in a small window around the given position.
- B. Lowercase character strings of tokens.
- C. Simple properties of characters inside tokens, such as capitalization, letters vs digits, punctuation, etc.
- D. Suffixes and prefixes of tokens with length 2 to 4 characters.

- E. Presence of tokens in local and global dictionaries, which contain words that were classified as certain entities someplace before – either anywhere (for global dictionaries), or in the current document (for local dictionaries).
- F. PoS tags of tokens.
- G. Stems of tokens.
- H. Presence of tokens in small manually prepared lists of semantic terms – such as months, days of the week, geographical features, company suffixes, etc.
- I. Presence of tokens inside gazetteers, which are huge lists of known entities.

The PoS tags are available only for the two CoNLL datasets, and the stems are available only for the CoNLL-D dataset. Both are automatically generated and thus contain many errors.

The gazetteers and lists of semantic terms are available for all datasets except CoNLL-D.

We tested the following feature sets:

- set0: checks properties A, B, C at the current and the previous token.
- set1: A, B, C, B+C in a window [-2...0].
- set2: A, B, C, B+C in a window [-2...+2].
- set2x: Same as set2, but only properties appearing > 3 times are used.
- set3: A, B, C, B+C in a window [-2...+2], D at the current token.
- set4: A, B, C, B+C in a window [-2...+2], D at the current token, E.
- set5: A, B, C, B+C, F, G in a window [-2...+2], D at the current token, E.
- set6: set4 or set5, H
- set7: set4 or set5, H, I

3.3 Hyperparameters

The MaxEntropy-based algorithms, MEMM and CRF, have similar hyperparameters, which define the priors for training the models. We experimented with two different priors – Laplacian (double exponential) $\text{Pr}_{\text{LAP}}(\boldsymbol{\lambda}) = \alpha \sum_i |\lambda_i|$ and Gaussian $\text{Pr}_{\text{GAU}}(\boldsymbol{\lambda}) = (\sum_i \lambda_i^2) / (2\sigma^2)$. Each prior depends upon a single hyperparameter specifying the “strength” of the prior. Note, that $\nabla \text{Pr}_{\text{LAP}}(\boldsymbol{\lambda})$ has discontinuities at zeroes of λ_i . Because of that, a special consideration must be given to the cases when λ_i approaches or is at zero. Namely,

- (1) if λ_i tries to change sign, set $\lambda_i := 0$, and allow it to change sign only on the next iteration, and
- (2) if $\lambda_i = 0$, and $\left| \frac{\partial}{\partial \lambda_i} L_T(\boldsymbol{\lambda}) \right| < \alpha$, do not allow λ_i to change, because it will immediately be driven back toward zero.

In some of the previous works (e.g., [22]) the Laplacian prior was reported to produce much worse performance than the Gaussian prior. Our experiments show them to perform similarly. The likely reason for the difference is poor handling of the zero discontinuities.

The RRM algorithm has three hyperparameters – the prior μ , the regularization parameter c , and the learning rate η .

4 Experimental Results

It is not possible to test every possible combination of algorithm, dataset and hyperparameter. Therefore, we tried to do a meaningful series of experiments, which would together highlight the different aspects of the algorithms.

All of the results are presented as final microaveraged F1 scores.

4.1 Experiment 1

In the first series of experiments we evaluated the dependence of the performance of the classifiers upon their hyperparameters. We compared the performance of the

Table 1. RRM results on CoNLL-E dataset

$\mu=0.01$	CoNLL-E_40K_set7			CoNLL-E_80K_set7			CoNLL-E_160K_set7		
	$c=0.001$	$c=0.01$	$c=0.1$	$c=0.001$	$c=0.01$	$c=0.1$	$c=0.001$	$c=0.01$	$c=0.1$
$\eta=0.001$	78.449	78.431	78.425	81.534	81.534	81.510	84.965	84.965	84.965
$\eta=0.01$	85.071	85.071	84.922	87.766	87.774	87.721	90.246	90.238	90.212
$\eta=0.1$	82.918	83.025	83.733	87.846	87.835	88.031	89.761	89.776	89.904
$\mu=0.1$									
$\eta=0.001$	84.534	84.552	84.534	87.281	87.281	87.264	89.556	89.556	89.573
$\eta=0.01$	85.782	85.800	85.800	89.032	89.032	89.066	91.175	91.175	91.150
$\eta=0.1$	82.439	82.709	83.065	63.032	63.032	63.032	30.741	30.741	56.445
$\mu=1.0$									
$\eta=0.001$	85.973	85.973	85.990	89.108	89.108	89.100	91.056	91.056	91.056
$\eta=0.01$	83.850	83.877	83.904	88.141	88.141	88.119	90.286	90.317	90.351
$\eta=0.1$	0	0	29.937	0	0	0	0	0	0

Table 2. RRM results on other datasets

$\mu=0.01$	CoNLL-D_20K_set7			MUC7_40K_set2x			CLF_80K_set2		
	$c=0.001$	$c=0.01$	$c=0.1$	$c=0.001$	$c=0.01$	$c=0.1$	$c=0.001$	$c=0.01$	$c=0.1$
$\eta=0.001$	43.490	43.490	43.453	48.722	48.722	48.650	49.229	49.229	49.244
$\eta=0.01$	46.440	46.438	46.472	63.220	63.207	62.915	64.000	64.040	63.710
$\eta=0.1$	44.878	44.943	45.995	61.824	62.128	63.678	58.088	58.628	61.548
$\mu=0.1$									
$\eta=0.001$	44.674	44.674	44.671	60.262	60.249	60.221	59.943	59.943	59.943
$\eta=0.01$	44.799	44.845	44.957	65.529	65.547	65.516	64.913	64.913	64.811
$\eta=0.1$	43.453	43.520	44.192	60.415	60.958	63.120	55.040	55.677	60.161
$\mu=1.0$									
$\eta=0.001$	44.682	44.682	44.694	66.231	66.231	66.174	65.408	65.408	65.408
$\eta=0.01$	43.065	43.080	43.195	62.622	62.579	62.825	59.197	59.311	59.687
$\eta=0.1$	0	0	6.123	2.922	2.922	8.725	0	0	1.909

Table 3. CRF results on a selection of datasets

CRF	CLF			CoNLL-D			MUC7	CoNLL-E
	20K_set2	40K_set2	80K_set2	40K_set1	80K_set1	160K_set1	80K_set0	80K_set0
GAU $\sigma = 1$	<u>76.646</u>	<u>78.085</u>	80.64	29.851	35.516	<u>39.248</u>	<u>80.756</u>	69.247
GAU $\sigma = 3$	75.222	77.553	79.821	28.530	35.771	38.254	80.355	<u>69.693</u>
GAU $\sigma = 5$	75.031	77.525	79.285	29.901	35.541	38.671	79.853	69.377
GAU $\sigma = 7$	74.463	77.633	79.454	30.975	<u>36.517</u>	38.748	79.585	69.341
GAU $\sigma = 10$	74.352	77.05	77.705	29.269	36.091	38.833	80.625	68.974
LAP $\alpha=0.01$	73.773	77.446	79.071	29.085	35.811	38.947	79.738	69.388
LAP $\alpha=0.03$	75.023	77.242	78.810	31.082	34.097	38.454	79.044	69.583
LAP $\alpha=0.05$	76.314	77.037	79.404	30.303	35.494	<u>39.248</u>	79.952	69.161
LAP $\alpha=0.07$	74.666	76.329	<u>80.841</u>	30.675	34.530	38.882	79.724	68.806
LAP $\alpha=0.1$	74.985	77.655	80.095	<u>31.161</u>	35.187	39.234	79.185	68.955

Table 4. MEMM results on a selection of datasets

MEMM	CLF			CoNLL-D			MUC7	CoNLL-E
	20K_set2	40K_set2	80K_set2	40K_set1	80K_set1	160K_set1	80K_set0	80K_set0
GAU $\sigma = 1$	<u>75.334</u>	<u>78.872</u>	<u>79.364</u>	<u>30.406</u>	35.013	<u>40.164</u>	<u>78.773</u>	67.537
GAU $\sigma = 3$	74.099	75.693	77.278	28.484	<u>35.330</u>	40.005	77.295	67.401
GAU $\sigma = 5$	73.959	74.685	77.316	28.526	35.043	39.799	77.489	67.870
GAU $\sigma = 7$	73.411	74.505	77.563	28.636	34.630	38.531	77.255	67.897
GAU $\sigma = 10$	73.351	74.398	77.379	28.488	33.955	37.830	77.094	<u>68.043</u>
LAP $\alpha=0.01$	71.225	74.04	75.721	28.316	34.329	40.074	78.312	67.871
LAP $\alpha=0.03$	72.603	72.967	76.540	29.086	35.159	38.621	77.385	67.401
LAP $\alpha=0.05$	71.921	75.523	75.370	30.425	33.942	39.984	78.262	67.908
LAP $\alpha=0.07$	72.019	74.486	77.197	30.118	35.250	39.195	76.646	67.833
LAP $\alpha=0.1$	72.695	75.311	76.335	30.315	33.487	40.861	78.141	67.421

classifiers on a selection of datasets, with different hyperparameter values. All of the algorithms showed moderate and rather irregular dependence upon their hyperparameters. However, single overall set of values can be selected.

The RRM results are shown in the Table 1 and the Table 2. As can be seen, selecting $\mu = 0.1$, $c = 0.01$ and $\eta = 0.01$ gives reasonably close to optimal performance on all datasets. All subsequent experiments were done with those hyperparameter values.

Likewise, the ME-based algorithms have no single best set of hyperparameter values, but have close enough near-optimal values. A selection of MEMM and CRF results is shown in the Table 3 and Table 4. For subsequent experiments we use CRF with Laplacian prior with $\alpha = 0.07$ and MEMM with Gaussian prior with $\sigma = 1$.

4.2 Training Size

In this series of experiments we evaluated the performance of the algorithms using progressively bigger training datasets: 10K, 200K, 400K, 800K and 1600K tokens. The results are summarized in the Fig.1. As expected, the algorithms exhibit very similar training size vs. performance behavior.

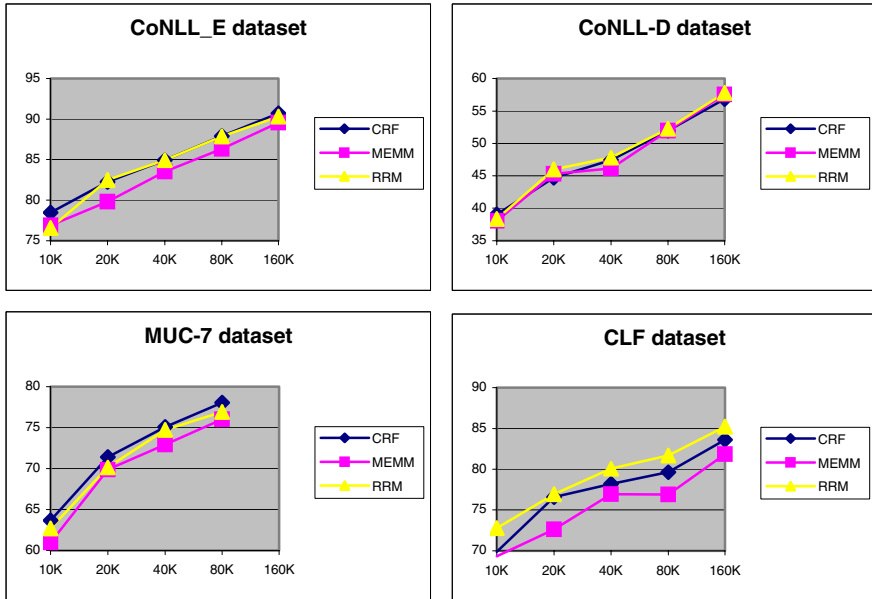


Fig. 1. Performance of the algorithms with different training sizes

Table 5. Performance of the algorithms with different feature sets

	MUC7			CoNLL-D			CoNLL-E		
	CRF	MEMM	RRM	CRF	MEMM	RRM	CRF	MEMM	RRM
set0	75.748	66.582	62.206	48.988	43.36	40.109	87.379	82.281	76.887
set1	75.544	67.075	68.405	50.672	49.164	48.046	87.357	82.516	81.788
set2	75.288	74.002	74.755	52.128	~52.01	51.537	86.891	87.089	87.763
set3	76.913	76.333	76.794	~60.172	59.526	61.103	88.927	88.711	89.110
set4	78.336	77.887	77.828	62.79	63.58	65.802	~90.037	~90.047	90.722
set5				~65.649	65.319	67.813	~90.139	~90.115	90.559
set6	78.969	78.442	78.016				~90.569	~90.492	90.982
set7	81.791	80.923	81.057				~91.414	90.88	91.777

4.3 Feature Sets

In this series of experiments we trained the algorithms with all available training data, but using different feature sets. The results are summarized in the Table 5. The results were tested for statistical significance using the McNemar test. All the perform-

ance differences between the successive feature sets are significant at least at the level $p=0.05$, except for the difference between set4 and set5 in CoNLL-E dataset for all models, and the differences between set0, set1, and set2 in CoNLL-E and MUC7 datasets for the CRF model. Those are statistically insignificant. The differences between the performance of different models that use same feature sets are also mostly significant. Exceptions are the numbers preceded by a tilde “~”. Those numbers are not significantly different from the best results in their corresponding rows.

As can be seen, both CRF and RRM generally outperform MEMM. Among the two, the winner appears to depend upon the dataset. Also, it is interesting to note that CRF always wins, and by a large margin, on feature sets 0 and 1, which are distinguished from the set 2 by absence of “forward-looking” features. Indeed, using “forward-looking” features produces little or no improvement for CRF, but very big improvement for local models, probably because such features help to alleviate the *label bias problem* [7].

5 Conclusions

We have presented the experiments comparing the three common state-of-the-art feature-rich probabilistic sentence classifiers inside a single system, using completely identical feature sets. The experiments show that both CRF and RRM significantly outperform MEMM, while themselves performing roughly similarly. Thus, it shows that the comparatively poor performance of CRF in the CoNLL 2003 NER task [16] is due to suboptimal feature selection and not to any inherent flaw in the algorithm itself.

Also, we demonstrated that the Laplacian prior performs just as well and sometimes better than Gaussian prior, contrary to the results of some of the previous researches.

On the other hand, the much simpler RRM classifier performed just as well as CRF and even outperformed it on some of the datasets. The reason of such surprisingly good performance invites further investigation.

References

1. Aitken, J. S.: Learning Information Extraction Rules: An Inductive Logic Programming approach. 15th European Conference on Artificial Intelligence. IOS Press. (2002)
2. Bikel, D. M., Schwartz, R., Weischedel, R.M.: An Algorithm that Learns What's in a Name. *Machine Learning*. (34): (1999) 211-231.
3. Chieu, H.L., Tou Ng, H.: Named Entity Recognition: A Maximum Entropy Approach Using Global Information. *Proceedings of the 17th International Conference on Computational Linguistics*. (2002)
4. McCallum, A., Freitag, D., Pereira, F.: Maximum Entropy Markov Models for Information Extraction and Segmentation. *Proceedings of the 17th International Conference on Machine Learning*. (2000)
5. Sun A., et al.: Using Support Vector Machine for Terrorism Information Extraction. 1st NSF/NIJ Symposium on Intelligence and Security Informatics. (2003)

6. Kushmerick, N., Johnston, E., McGuinness, S.: Information extraction by text classification. IJCAI-01 Workshop on Adaptive Text Extraction and Mining. Seattle, WA. (2001)
7. Lafferty, J., McCallum, A., Pereira, F.: Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. Proc. 18th International Conf. on Machine Learning. (2001)
8. Rosenfeld, B., Feldman, R., Fresko, M., Schler, J., Aumann, Y.: TEG - A Hybrid Approach to Information Extraction. Proc. of the 13th ACM. (2004)
9. Berger, A., della Pietra, S., della Pietra, V.: A maximum entropy approach to natural language processing. *Computational Linguistics* 22(1), (1996) 39-71.
10. Darroch, J.N., Ratcliff, D.: Generalized iterative scaling for log-linear models. *Annals of Mathematical Statistics* (1972). 43(5): 1470-1480.
11. Borthwick, A., Sterling, J., Agichtein, E., Grishman, R.: Exploiting Diverse Knowledge Sources via Maximum Entropy in Named Entity Recognition. In the proceedings of the 6th Workshop on Very Large Corpora. (1998)
12. Kim Sang, T., Erik, F., De Meulder, F.: Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition. Edmonton, Canada (2003)
13. Florian, R., Ittycheriah, A., Jing, H., Zhang, T.: Named Entity Recognition through Classifier Combination. In: Proceedings of CoNLL-2003, Edmonton, Canada, (2003), pp. 168-171.
14. Zhang, T., Johnson, D.: A Robust Risk Minimization based Named Entity Recognition System. In: Proceedings of CoNLL-2003, Edmonton, Canada, (2003), pp. 204-207.
15. Chieu, H.L., Tou Ng, H.: Named Entity Recognition with a Maximum Entropy Approach. In: Proceedings of CoNLL-2003, Edmonton, Canada, (2003), pp. 160-163.
16. McCallum, A., Li, W.: Early results for Named Entity Recognition with Conditional Random Fields, Feature Induction and Web-Enhanced Lexicons. In: Proceedings of CoNLL-2003, Edmonton, Canada, (2003), pp. 188-191.
17. Joachims, T.: Learning to Classify Text Using Support Vector Machines. Dissertation, Kluwer, (2002)
18. Nigam, K., Lafferty, J., McCallum, A.: Using maximum entropy for text classification. In IJCAI-99 Workshop on Machine Learning for Information (1999) 61—67.
19. Sha, F., Pereira, F.: Shallow parsing with conditional random fields, Technical Report CIS TR MS-CIS-02-35, University of Pennsylvania, (2003)
20. Zhang, T., Damerau, F., Johnson, D.: Text Chunking using Regularized Winnow. Meeting of the Association for Computational Linguistics. (2001) 539-546
21. Zhang, T., Regularized Winnow Methods. NIPS, (2000) 703-709
22. Peng, F., McCallum, A.: Accurate Information Extraction from Research Papers Using Conditional Random Fields. (1997)
23. Chinchor, N. MUC-7 Named Entity Task Definition Dry Run Version, Version 3.5. Proceedings of the Seventh Message Understanding Conference. (1998)
24. Borthwick, A., Sterling, J., Agichtein, E., Grishman, R. Description of the MENE Named Entity System as Used in MUC-7. Proceedings of the Seventh Message Understanding Conference. (1998)

Knowledge Discovery from User Preferences in Conversational Recommendation

Maria Saa..., James Reilly, Lorraine McGinty, and Barry Smyth

Smart Media Institute, University College Dublin,
Belfield, Dublin 4, Ireland

{maria, james.d.reilly, lorraine.mcginity, barry.smyth}@ucd.ie

Abstract. Knowledge discovery for personalizing the product recommendation task is a major focus of research in the area of conversational recommender systems to increase efficiency and effectiveness. Conversational recommender systems guide users through a product space, alternatively making product suggestions and eliciting user feedback. Critiquing is a common and powerful form of feedback, where a user can express her feature preferences by applying a series of directional critiques over recommendations, instead of providing specific value preferences. For example, a user might ask for a ‘less expensive’ vacation in a travel recommender; thus ‘less expensive’ is a critique over the *price* feature. The expectation is that on each cycle, the system discovers more about the user’s *soft* product preferences from minimal information input. In this paper we describe three different strategies for knowledge discovery from user preferences that improve recommendation efficiency in a conversational system using critiquing. Moreover, we will demonstrate that while the strategies work well separately, their combined effort has the potential to considerably increase recommendation efficiency even further.

1 Introduction

Recommender systems are a key component of knowledge discovery. Each time a user decides which product to purchase, she is faced with a large number of choices. In this paper, we focus on the problem of recommending products to users based on their preferences. We describe a system that uses a knowledge discovery process to learn about user preferences from a series of critiques. A critique is a directional critique over a recommendation, where the user expresses her preference for a feature. For example, a user might say ‘less expensive’ to indicate that she prefers a product that is less expensive. The system then uses this information to refine its recommendations. We describe three different strategies for knowledge discovery from user preferences that improve recommendation efficiency in a conversational system using critiquing. Moreover, we will demonstrate that while the strategies work well separately, their combined effort has the potential to considerably increase recommendation efficiency even further.

Recalling the decision capabilities of the feedback has been shown; e.g. the decision capabilities of the feedback [7]. In this case, the decision capabilities of the feedback can be used [8], the decision capabilities of the feedback can be used [8], the decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8].

With the decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8].

In this case, the decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8].

2 Background

This section describes the decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8].

The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8]. The decision capabilities of the feedback can be used [8].


```

q: query, CB: CaseBase, cq: critique, cr : current recommendation, U : User model

1.  define Incremental_Critiquing(q, CB)
2.  cq:= null
3.  U:= null
4.  begin
5.  do
6.  cr ← ItemRecommend(q, CB, cq, U)
7.  cq ← UserReview(cr, CB)
8.  q ← QueryRevise(q, cr)
9.  U ← UpdateModel(U, cq, cr)
10. until UserAccepts(cr)
11. end

12. define ItemRecommend(q, CB, cq, U)
13. CB' ← {c ∈ CB | Satisfies(c, cq)}
14. CB' ← sort cases in CB' in decreasing Quality
15. cr ← most quality case in CB'
16. return cr

17. define UserReview(cr, CB)
18. cq ← user critique for some f ∈ cr
19. CB ← CB - cr
20. return cq

21. define QueryRevise(q, cr)
22. q ← cr
23. return q

24. define UpdateModel(U, cq, cr)
25. U ← U - contradict(U, cq, cr)
26. U ← U - refine(U, cq, cr)
27. U ← U + (<cq, cr>)
28. return U
    
```

Fig. 1. The incremental critiquing algorithm

he rec. e. da. l. e. l. , . . . c. l. l. e. i. a. a. e. a. . . i. e. c. e. h. e. e. c. c. e. A. l. l. i. e. d. e. l. f. h. e. i. c. e. e. a. c. l. l. i. g. a. g. i. h. i. g. i. e. i. Fig. 1.

The i. c. e. e. a. c. l. l. i. g. a. g. i. h. c. l. l. i. f. 4. e. e. e. : (1) a. l. e. c. a. e. c_r i. . . . h. e. e. b. a. e. d. . . h. e. c. . . e. . . a. d. h. e. . . e. i. . . c. l. l. e. ; (2) h. e. . . . h. e. e. c. . . e. d. a. l. . . a. d. a. n. e. a. d. i. e. c. i. a. f. e. a. c. l. l. e. , *cq*; (3) h. e. . . , *q* i. . . . f. . . h. e. e. c. c. e.; (4) h. e. . . d. e. , *U* i. . . d. a. e. d. b. a. d. d. i. g. h. e. a. c. l. l. i. e. *cq* a. d. . . i. g. a. h. e. c. l. l. i. e. h. a. a. e. i. c. . . i. e. . . i. h. i. . The . . . e. d. a. l. . . c. e. . . i. a. e. e. h. e. h. e. h. e. . . e. . . e. e. d. i. h. a. i. a. b. e. c. a. e. . . h. e. h. e. g. r. e. . .

I. . . . a. . . , h. e. e. c. . . e. d. a. l. . . c. e. l. l. i. . e. e. c. d. b. h. e. . . e. . . d. e. l. f. . . e. i. . . c. l. l. e. , *U*, h. a. i. i. c. e. e. a. . . d. a. e. d. . . e. a. c. c. e. I. c. e. e. a. c. l. l. i. g. . . d. i. e. h. e. b. a. i. c. c. l. l. i. g. a. g. i. h. . . I. e. a. d. f. . . d. e. i. g. h. e. . . e. d. c. a. e. . . h. e. b. a. i. f. h. e. i. . . i. a. i. . . h. e. e. c. . . e. d. c. a. e. i. a. . . c. . . e. a. , c. c. e. (e. e. E. a. i. . . 1) f. . . e. a. c. i. d. a. e. c. a. e. . The . . . a. i. b. i. l. . . c. e. i. e. e. i. a. . . h. e. e. c. e. a. g. e. f. c. l. l. i. e. i. h. e. . . e. . . d. e. h. a. h. i. c. a. e. . a. i. e. . The . . . h. e. c. . . a. i. b. i. l. . . c. e. a. d. h. e. c. a. d. i. d. a. e'. (c') i. . . i. a. i. . . h. e. c. . . e. . . e. c. . . e. d. a. l. . . (c_r) a. e. c. . . b. i. e. d. i. . . d. e. . . b. a. . . a. . . e. a. c. c. e. (e. e. E. a. i. . . 2, b. d. e. f. a. . . β = 0.75). The . . . a. i. . . c. e. i. e. d. . . a. . . h. e. . . e. e. d. c. a. e. . . i. . . h. e. e. . . e. c. . . e. d. a. l. . . c. c. e. (e. e. i. e. 14.1. Fig. 1) a. d. h. e. c. a. e. i. h. h. e. h. i. g. h. e. . . a. i. . . h. e. c. h. e. a. h. e. e. . . e. c. . . e. d. a. l. . .

$$Compatibility(c', U) = \frac{\sum_{U_i} satisfies(U_i, c')}{|U|} \tag{1}$$

$$Quality(c', c_r, U) = \beta \cdot Compatibility(c', U) + (1 - \beta) \cdot Similarity(c', c_r) \tag{2}$$

A. g. i. h. i. l. a. i. a. l. a. c. l. l. i. e. b. a. e. d. . . e. . . d. e. h. i. c. h. i. c. . . e. d. f. h. e. c. l. l. i. e. h. a. h. a. e. b. e. e. c. h. e. b. h. e. . . e. . . f. a. . . O. e. f. h. e. . . i. . .

finded. Hence, the accuracy of each approach is affected by the number of high-level features. Hence, we compare the accuracy of each approach. The results are as follows: (1) the accuracy of each approach is affected by the number of features; (2) the accuracy of each approach is affected by the number of features.

3 Knowledge Discovery Strategies

The accuracy of each approach is affected by the number of features. Hence, we compare the accuracy of each approach. The results are as follows: (1) the accuracy of each approach is affected by the number of features; (2) the accuracy of each approach is affected by the number of features.

3.1 Discovering Satisfactory Cases: Highest Compatibility Selection

As we have seen, the accuracy of each approach is affected by the number of features. Hence, we compare the accuracy of each approach. The results are as follows: (1) the accuracy of each approach is affected by the number of features; (2) the accuracy of each approach is affected by the number of features.

Figure 2 shows the accuracy of each approach. As we have seen, the accuracy of each approach is affected by the number of features. Hence, we compare the accuracy of each approach. The results are as follows: (1) the accuracy of each approach is affected by the number of features; (2) the accuracy of each approach is affected by the number of features.

The compatibility function. We have concluded that the accuracy of each approach is affected by the number of features. Hence, we compare the accuracy of each approach. The results are as follows: (1) the accuracy of each approach is affected by the number of features; (2) the accuracy of each approach is affected by the number of features.

```

q: query, CB: CaseBase, cq: critique, cr: current recommendation, U: User Model

1.define ItemRecommend(q, CB, cq, U)
2. CB' ← {c ∈ CB | Satisfies(c, cq)}
3. CB' ← sort cases in CB' in decreasing compatibility score
4. CB'' ← selects those cases in CB' with highest compatibility
5. CB'' ← sort cases in CB'' in decreasing order of their sim to q
6. cr ← most similar case in CB''
7.return cr
    
```

Fig. 2. Adapting the incremental critiquing algorithm *ItemRecommend* procedure to improve focus on recommendation by using Highest Compatibility Selection strategy

... a i f h e ... e f e e c e . F ... h i e a ... , e e a a e h e e a i g c a e a i f h e ... e e a e f a e i a ... (RLP) [12], h i c h c ... i ... f a i i g h e ... f f ... e e a d i a e f a e . R e f f e r e n c e e . L e a ... i g h e ... i ... a b a e d ... (FMDP).

Each case is evaluated according to a binary criterion based on each case's rating. The evaluation function (see Example 3). This function evaluates the ... of each case ... the ... the ... e e f e a i g c a e e a ... e h a c e ... a c c i d i g ... h e c i i e h e e h a ... e e c e d .

$$\text{Compatibility}(c', U_f) = \begin{cases} \text{comp}(c') + \alpha \times (1 - \text{comp}(c')) & \text{if } \dots U_f \\ \text{comp}(c') + \alpha \times (0 - \text{comp}(c')) & \text{if } \dots U_f \end{cases} \quad (3)$$

O ... g a i ... a i a ... a i f a h e ... e f e e c e . T h ... e a e ... i g f ... a e f a i a c ... a i b e c a e (i.e., h e c a e ... h i c h h a e h e h i g h e c ... a i b i i (... ,) a e c ... i d e i g a h e ... e f e e c e (U) ... a c i ... i e) . A ... h e b e g i ... i g f e a c h e ... e a c h c a d i d a e c a e , ' h a a d e f a c ... a i b i i ... a e (i.e., ... , (') ...) . T h ... a e i ... d a e d ... e c c e a - i g i ... a c c ... h e a i f a c i ... f h e c ... e c i i e . T h e ... a a e e , i E a i . 3 i h e e a ... i g a e h i c h ... a ... e ... 0.1 , 0.2 a e ; a a g e ... a e e a d ... a a g e g a b e e e c a e i e a ... a g e . I ... c a e , h e e a ... i g a e i ... i ... a ... i c e e a e ... i g f ... e e f a i f a c i ... I ... h e ... d , e a e ... , i g ... b a i a e f a e h a a ... i e a ... i c a ... i b e ... a 1.0 a e , a ... a ... i d e i RLP.

I ... i ... a ... e h a E a i . 3 d a e h e c ... a i b i i ... a e ... e d b e a c h c a e a c c i d i g ... h e a ... e c i i e (U_f) a ... e d ... c ... i g a h e e f c i i e i e h e i c e e a a ... a c h (e e E a i . 1) . T h e *Compatibility*(c', U_f) a e c ... e d i h e c ... e c c e i b e h e (comp(c')) i h e e c c e .

3.2 Discovering Important Features: Local User Preference Weighting

The ... e i ... a e g h i g h ... h e c a e d i e ... i a i ... b e . I ... h e ... d , i i f c e d ... d i c e i g c a e h a ... a i a ... a i f ... e ... e f e e c e .

Now, we use the following heuristic to compute the feature f importance. We use the following heuristic, (LW), a feature f is important if each feature f is important in each case c a high importance factor. If the feature f is important in all cases, we use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c .

Overall, we use the following heuristic to compute the feature f importance. We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c .

$$weight(c'_f) = 1 - \left(\frac{\#critiques\ in\ U\ that\ satisfy\ feature\ f\ in\ case\ c'}{\#total\ critiques\ feature\ f\ in\ U} \times 0.5 \right) \tag{4}$$

We get a feature f importance of each case c . The feature f importance of each case c is the average of the importance of each case c . The feature f importance of each case c is the average of the importance of each case c . The feature f importance of each case c is the average of the importance of each case c . The feature f importance of each case c is the average of the importance of each case c . The feature f importance of each case c is the average of the importance of each case c .

Overall, we use the following heuristic to compute the feature f importance. We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c .

$$Similarity(c', c_r) = \sum_{\forall_f} weight(c'_f) \times similarity(c'_f, c_{r_f}) \tag{5}$$

The similarity between the candidate case c' and the reference case c_r for each feature f is computed as the average of the importance of each feature f . The similarity between the candidate case c' and the reference case c_r is the average of the similarity between the candidate case c' and the reference case c_r for each feature f .

3.3 Discovering Query Knowledge: Binary Search

The first step in the binary search algorithm is to find the most important features. We use the following heuristic to compute the feature importance. We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c . We use the following heuristic, (LW), a feature f is important if each case c is important in each case c .

ab, e he, edia). We ace he c, 11 ed fea , e 1 a ce di g a e , de a d
 d he idde a e if he be, f ca e 1 dd , d he idde a i a d
 c. e he ea a a e be ee he if e ha e a e e . be, f ca e .

O e 1 , a 1 , ha a . eed , be c, ide, ed, 1 e 1 c, 11 e .
 he a e fea , e. F, e a e , e ha a e ha a ed i a , e 1 c ce
 f, a , , aca 1 . ha a 2500, ec. e da 1 a d, 1 he c, e c -
 ce, he e a ha he efe, a , ha a 1000 aca 1 . I ch
 1 a 1 1 he c, e c ce, a he ca e i c di g h e ha e eed a 2500
 aca 1 . 1 a i f he c, e c, 11 e , ha 1000. If ec -
 e he edia a e age 1 he ea ch ace, e a i c de h e
 ca e e ec ed e 1 . b he e. T a id he e 1 a 1 , e e he h
 f c, 11 e a i ed b he e 1 , de c c, ec . he ea ch ace. The
 e 1 c, 11 e ed 1 he e . de a e ea ed a a e f
 [13] ha a c . he be, f e a i g ca e ha 1 be ed .
 c. e he edia a e .

S, f . i g he ea , e a e , e . c, ide, c . i g he edia f
 h e ca e ha a e ha 1000 a d ha 2500.
 A de a ed 1 i e 5 f Fig e 3, be, f e c . i g he edia , e chec f, he
 e 1 e ce, f e 1 . a i ed c, 11 e ha c e he i c 1 . f ca e i CB',
 a de 1 i a e he e ca e f . f, he c, ide, a 1 . P di e e , e e 1 ,
 c, 11 e . decide ha ca e h d be c, e ed b CB', a d . 1 a e e he
 a e e ec 1 . b . d f, f_{eq}.

The e . 1 a 1 . behi d e e 1 . i ce e a c i -
 1 i g a . ed ce c, 11 e e e 11 e e ce, a d i . e, ec. e da 1 .
 e cie c b di c, e i g a i fa c d c, f . e e a id . I h , h
 bi a . ea ch . ea . ache a be he, ec. e de f c, 11 e a ch . h e
 ca dida e ca e ha : (1) a i f he c, e c, 11 e ; (2) f . e 1 . a i ed
 c, 11 e ; a d (3) a e 1 i a . he c, e ca e b f, he a a f . 1 , a d h .
 ha e he ca a b 1 . f a i g a i g he ea ch ace f . 1 c .

4 Evaluation

I h i a e . fa e ha e a g ed ha he i ce e a f f c, 11 i g 1
 1 i ed b 1 . e de c . ec. e d ca e ha d a i a . a i f he
 . e . efe, e ce . We e h ee . a e g i e ha a id edge di c e 1 a
 e . 1 e, e ne a acc , ac a d, ec. e da 1 . e cie c . Th e c 1 .
 de c, i be he ea ed e a a 1 . e h d i g ha e ed a d he e . ha
 e . ed .

4.1 Setup

The e a a 1 . a e f . ed i g he a da d T a e da a e (a a i a be f .
 ,) h i ch c . 1 . f 1024 aca 1 . ca e . Each ca e 1 de -
 c, i bed i . e . f 9 fea , e i c di g , , ec . The da a e a
 ch e be ca e 1 c . a i e i ca a d i a fea , e a d i a ide
 a i de ea ch ace .

We evaluate the highest capability (HCS), the calculation efficiency (CE), the load (LW), the bandwidth (BS) and the average capability (ALL) of the device (ALL) of the device (ACEE) (ACEE).

4.2 Methodology

We divided the hardware into three parts: the hardware, the software, and the hardware. According to each case (which was called the 'base') the hardware is divided into three parts. First, the hardware is divided into three parts: the hardware, the hardware, and the hardware. Second, the hardware is divided into three parts: the hardware, the hardware, and the hardware. Third, the hardware is divided into three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware.

Regarding the hardware [15] has high speed of calculation and high efficiency [14] and the hardware is divided into three parts: the hardware, the hardware, and the hardware. The hardware is divided into three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware.

4.3 Recommendation Efficiency

We evaluate the hardware efficiency (CE) which is the average efficiency of the hardware (ALL) of the device (ACEE) (ACEE). The hardware is divided into three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware. Each generation is a set of three parts: the hardware, the hardware, and the hardware.

Figure 4(B) highlights the benefit of each strategy (HCS, LW, and BS) relative to the combined strategy (ALL) in terms of reducing the cycles required to solve the incremental problem. We find that the ALL strategy reduces the cycles by 2.65% and 7.5%, which are similar to the relative benefits of the HCS, LW and BS approaches. The relative benefit of the highest performing technique (HCS) approaches, which range between 2.65% and 3.81%, because it decreases the average number of cycles required to solve each problem. However, the ALL strategy is significantly better than each of the individual strategies (HCS, LW, and BS) on the whole. On the whole, the combined approach (LW) gives the highest benefit, ranging from 4.5% to 6.73%, the average value. The overall benefit of the ALL strategy is significantly better than the individual strategies and decreases the average number of cycles.

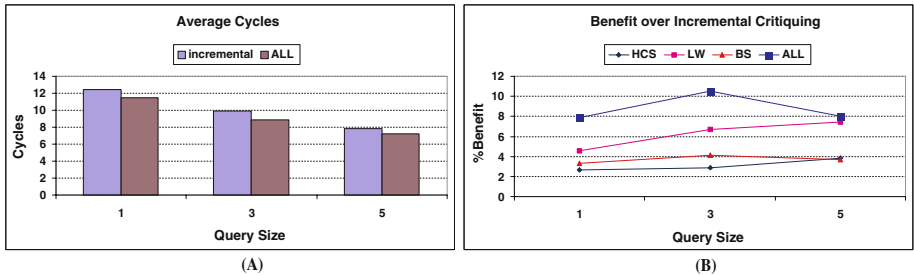


Fig. 4. Average session length and benefit over incremental critiquing

On the whole, the combined strategy is significantly better than the individual strategies. The ALL strategy is significantly better than each of the individual strategies (HCS, LW, and BS) on the whole. On the whole, the combined approach (LW) gives the highest benefit, ranging from 4.5% to 6.73%, the average value. The overall benefit of the ALL strategy is significantly better than the individual strategies and decreases the average number of cycles.

In terms of highlighting the benefit of the combined strategy relative to the individual strategies, see Figure 5. We have observed that the ALL strategy is significantly better than each of the individual strategies (HCS, LW, and BS) on the whole. On the whole, the combined approach (LW) gives the highest benefit, ranging from 4.5% to 6.73%, the average value. The overall benefit of the ALL strategy is significantly better than the individual strategies and decreases the average number of cycles.

(e.g., 84%) e.e. i.h.h., e.e. i.h.h. he.e he BS a... ach d e... ha e... ch f.a... ,... i... a ec... he, ec... e da... .

T... a i.e, a i.g i ca... e cie c be e... i e... ed b HCS, LW a d BS... a e.gie, he c... a ed... he i.c.e e... a c, i i g a... ach. The... a i c... i b... i f h i... a e... i ha... he... ed... a e.gie a... i... he d i.c... e... f... e f... ec... e da... . edge, a... i g he... e... , i, i i e... d c... ha b e... a i f... he... e. We ha e de... , a ed ha... h i... a... ach i high e ec i.e, e e... i... i a... . he e... a... i i a... edge f... e... efe, e ce... i a a i a b e (e.g., c, i i i g a... ach). F... he... e, he... e... f... he c... - b i e d... a e.gie... h... a i g i ca... i c, e a e i... ec... e da... e cie c... he c... a ed... i c, e e... a c, i i i g a d a... . he b a i c c, i i i g a... ach... ed b [11].

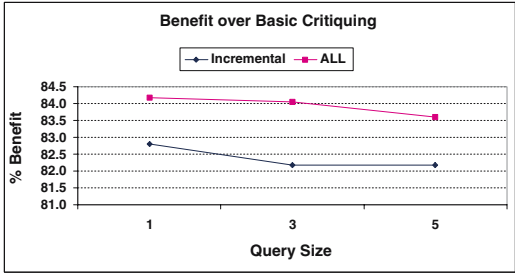


Fig. 5. Incremental and ALL benefit over basic critiquing approach

5 Conclusions

The d i.c... e... f... i... c i... e... , e f... e... c... e... edge... i... e... e... a... . decide... h i c... d... c... e... c... e... e... da... . a... e... h... e... i... a... b e f... , each... e... d... i... g... a... i... e... c... -... e... i... e... a... c... . I... h... i... a... e... e... h... a... e... ed... h... e... d... i... c... e... ,... a... e.gie... h... a... i... a... i... ,... e... e... c... e... e... da... . e... cie... c... . F... i... f... a... ,... e... h... a... e... e... ed... a... c... a... e... i... i... i... a... i... ,... a... e.g... h... a... a... i... i... e... h... e... e... ,... e f... e... c... e... e... i... e... . S... e... c... o... d... ,... e... h... a... e... e... ed... a... e... ,... e f... e... c... e... e... i... g... i... g... a... e.g... h... a... i... i... i... e... f... e... a... . c... a... . each... c... a... e... . F... i... a... ,... e... h... a... e... e... ed... a... e... ,... a... e.g... h... a... h... a... h... e... c... a... b... i... i... f... . a... i... g... h... e... a... c... h... . a... c... .

O... e... e... i... e... i... d... i... c... a... h... a... h... e... h... e... e... a... h... a... h... e... e... i... a... d... e... i... e... ,... h... h... i... e... c... i... c... b... e... . R... e... d... c... i... . i... h... e... a... e... a... g... e... g... h... f... ,... e... c... e... e... d... a... i... . e... i... i... e... e... e... d... i... a... f... h... e... ,... a... ,... b... h... e... a... a... e... a... d... c... o... b... i... e... d... ,... h... e... c... . a... e... d... . h... e... i... c... e... e... a... a... d... b... a... i... c... ,... i... i... g... e... . I... ,... a... ,... h... e... ,... e... d... . a... e.gie... a... e... . c... i... e... ,... g... e... a... . b... e... a... i... c... a... b... e... a... . a... i... d... e... a... g... e... f... e... c... e... e... da... . i... c... e... a... i... . I... a... i... c... a... ,... h... e... h... a... a... . e... a... c... -... e... e... d... c... -... a... c... h... e... e... ,... e... c... e... e... da... . i... e... i... i... a... e... i... e... . b... e... ,... a... c... e... d... ,... a... d... /... d... . a... i... . h... e... e... i... i... a... e... ,... f... e... e... b... a... c... i... i... e... . b... e... a... i... a... b... e... .

References

1. D.W. Aha, L.A. Breslow, and H. Muñoz-Avila. Conversational Case-Based Reasoning. *Applied Intelligence*, 14:9–32, 2000.
2. D. McSherry. Increasing Dialogue Efficiency in Case-Based Reasoning without Loss of Solution Quality. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence*, pages 121–126. Morgan-Kaufmann, 2003.
3. D. McSherry and C. Stretch. Automating the Discovery of Recommendation Knowledge. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence*, page forthcoming. Morgan-Kaufmann, 2005.
4. R. Burke, K. Hammond, and B.C. Young. The FindMe Approach to Assisted Browsing. *Journal of IEEE Expert*, 12(4):32–40, 1997.
5. L. McGinty and B. Smyth. Comparison-Based Recommendation. In Susan Crow, editor, *Proceedings of the 6th European Conference on Case-Based Reasoning*, pages 575–589. Springer, 2002. Aberdeen, Scotland.
6. H. Shimazu. ExpertClerk: A Conversational Case-Based Reasoning Tool for Developing Salesclerk Agents in E-Commerce Webshops. *Artificial Intelligence Review*, 18(3-4):223–244, 2002.
7. B. Smyth and L. McGinty. An Analysis of Feedback Strategies in Conversational Recommender Systems. In P. Cunningham, editor, *Proceedings of the 14th National Conference on Artificial Intelligence and Cognitive Science*, 2003. Dublin, Ireland.
8. L. McGinty and B. Smyth. Tweaking Critiquing. In *Proceedings of the Workshop on Personalization and Web Techniques at the International Joint Conference on Artificial Intelligence*. Morgan-Kaufmann, 2003.
9. R. Burke. Interactive Critiquing for Catalog Navigation in E-Commerce. *Artificial Intelligence Review*, 18(3-4):245–267, 2002.
10. J. Reilly, K. McCarthy, L. McGinty, and B. Smyth. Incremental Critiquing. In M. Bramer, F. Coenen, and T. Allen, editors, *Research and Development in Intelligent Systems XXI. Proceedings of AI-2004*, pages 101–114. Springer, 2004. Cambridge, UK.
11. R. Burke, K. Hammond, and B. Young. Knowledge-Based Navigation of Complex Information Spaces. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 462–468. AAAI Press/MIT Press, 1996. Portland, OR.
12. M.E. Harmon and S.S. Harmon. Reinforcement learning: A tutorial, 1996.
13. M. Stolze. Soft Navigation in Electronic Product Catalogs. *International Journal on Digital Libraries*, 3(1):60–66, 2000.
14. B. Smyth and L. McGinty. The Power of Suggestion. In *Proceedings of the International Joint Conference on Artificial Intelligence*. Morgan-Kaufmann, 2003.
15. K. McCarthy, L. McGinty, B. Smyth, and J. Reilly. On the Evaluation of Dynamic Critiquing: A Large-Scale User Study. In *Proceedings Twentieth National Conference on Artificial Intelligence*, pages 535–540. AAAI Press / The MIT Press, 2005.

Unsupervised Discretization Using Tree-Based Density Estimation

Gabi Schmidberger and Eibe F. F. de Azevedo

Department of Computer Science, University of Waikato,
Hamilton, New Zealand
{gabi, eibe}@cs.waikato.ac.nz

Abstract. This paper presents an unsupervised discretization method that performs density estimation for univariate data. The subintervals that the discretization produces can be used as the bins of a histogram. Histograms are a very simple and broadly understood means for displaying data, and our method automatically adapts bin widths to the data. It uses the log-likelihood as the scoring function to select cut points and the cross-validated log-likelihood to select the number of intervals. We compare this method with equal-width discretization where we also select the number of bins using the cross-validated log-likelihood and with equal-frequency discretization.

1 Introduction

Discretization is a well-known technique for data analysis. It is a simple and effective way to reduce the dimensionality of data. The resulting discrete data can be used for a variety of tasks, such as classification, clustering, and visualization. The most common discretization methods are equal-width, equal-frequency, and entropy-based. Each method has its own advantages and disadvantages. In this paper, we propose a new method for discretization based on tree-based density estimation. The idea is to use a decision tree to estimate the density of the data. The resulting tree structure can be used to define the bins of a histogram. The number of bins is determined by the number of internal nodes in the tree. The width of each bin is determined by the width of the corresponding leaf node. This method is simple and effective, and it can be used for a wide range of data sets.

Our goal is to provide a simple and effective method for discretization. The resulting discrete data can be used for a variety of tasks, such as classification, clustering, and visualization. The most common discretization methods are equal-width, equal-frequency, and entropy-based. Each method has its own advantages and disadvantages. In this paper, we propose a new method for discretization based on tree-based density estimation. The idea is to use a decision tree to estimate the density of the data. The resulting tree structure can be used to define the bins of a histogram. The number of bins is determined by the number of internal nodes in the tree. The width of each bin is determined by the width of the corresponding leaf node. This method is simple and effective, and it can be used for a wide range of data sets.

$$h = \frac{n_i}{w_i * N} \quad (1)$$

The algorithm is defined as follows. In Section 2, we discuss the basic definitions and the algorithm. In Section 3, we describe the algorithm. In Section 4, we describe the algorithm. In Section 5, we describe the algorithm. In Section 6, we describe the algorithm. In Section 7, we describe the algorithm.

2 Non-parametric Density Estimation

Definition 1. Let X be a random variable with density function f . Let X_1, \dots, X_n be a random sample of size n from X . Let \hat{f}_n be the kernel density estimator of f with kernel K and bandwidth h . Let μ and σ^2 be the mean and variance of X . The kernel density estimator \hat{f}_n is unbiased for f if and only if $\int K(u) du = 1$ and $\int u^2 K(u) du = 0$. The kernel density estimator \hat{f}_n is consistent for f if $\int K(u) du = 1$ and $\int u^2 K(u) du = 0$.

Non-parametric density estimation is a branch of statistics that deals with the estimation of the probability density function of a random variable without making any assumptions about the form of the distribution. The most common method for non-parametric density estimation is the kernel density estimator. The kernel density estimator is a non-parametric estimator of the probability density function of a random variable. It is based on the idea of smoothing the empirical distribution function. The kernel density estimator is defined as follows. Let X_1, \dots, X_n be a random sample of size n from a random variable X with density function f . Let K be a kernel function and h be the bandwidth. The kernel density estimator \hat{f}_n is defined as follows:

The kernel density estimator is defined as follows. Let X_1, \dots, X_n be a random sample of size n from a random variable X with density function f . Let K be a kernel function and h be the bandwidth. The kernel density estimator \hat{f}_n is defined as follows:

Let X_1, \dots, X_n be a random sample of size n from a random variable X with density function f . Let K be a kernel function and h be the bandwidth. The kernel density estimator \hat{f}_n is defined as follows:

3 Existing Unsupervised Discretization Methods

We consider the standard histogram-based discretization method: divide the range of the variable into bins and assign each data point to the bin it falls into. The resulting histogram is used to estimate the probability density function of the variable. A common choice for the histogram is the equal-width histogram. The resulting histogram is used to estimate the probability density function of the variable. We use [4].

For the histogram-based method, the histogram is used to estimate the probability density function of the variable. The histogram is used to estimate the probability density function of the variable. We use [4].

The histogram-based method is used to estimate the probability density function of the variable. The histogram is used to estimate the probability density function of the variable. We use [4].

4 Cross-Validating the Log-Likelihood

Cross-validation is used to estimate the probability density function of the variable. The histogram is used to estimate the probability density function of the variable. We use [4].

The histogram-based method is used to estimate the probability density function of the variable. The histogram is used to estimate the probability density function of the variable. We use [4].

Let n_i be the number of data points in bin i , n_{i-test} be the number of data points in bin i used for testing, w_i be the width of bin i , and N be the total number of data points. The histogram-based log-likelihood is given by:

$$L = \sum_i n_{i-test} * \log \frac{n_i}{w_i * N} \tag{2}$$

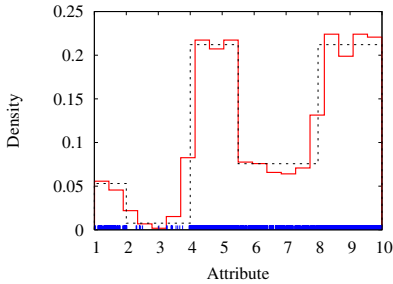


Fig. 1. Equal-width method with 20 bins

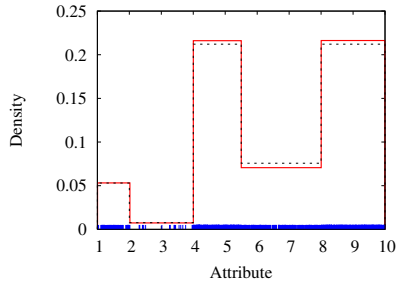


Fig. 2. Our TUBE-Method chose 5 bins of varying length

The error between the bins. If n_i is the number of data points in the bin i , then the error is the sum of the absolute differences between the density estimate and the true density for each bin. The error is defined as $L = \sum_i |n_i - test * \log \frac{n_i + \frac{w_i}{W}}{w_i * (N + 1)}|$.

$$L = \sum_i n_{i-test} * \log \frac{n_i + \frac{w_i}{W}}{w_i * (N + 1)} \quad (3)$$

5 Tree-Based Unsupervised Discretization

The proposed method is based on the idea of using a decision tree to discretize the data. The tree is trained on the data, and the resulting tree structure is used to discretize the data. The tree structure is used to discretize the data into bins of varying length. The proposed method is based on the idea of using a decision tree to discretize the data. The tree is trained on the data, and the resulting tree structure is used to discretize the data. The tree structure is used to discretize the data into bins of varying length.

We call the proposed method TUBE (Tree-based Unsupervised Bin Estimation), because it is a tree-based method for unsupervised bin estimation. The method is based on the idea of using a decision tree to discretize the data. The tree is trained on the data, and the resulting tree structure is used to discretize the data. The tree structure is used to discretize the data into bins of varying length.


```

maxNumBins           = [find optimal number of splits];
numSplits            = 0;
splitPriorityQueue    = empty;

firstBin             = new Bin(bin that contains the whole attribute range);
fringe               = [initialize with (firstBin)];

REPEAT {
  FOR (bin = all bins in fringe) {
    split = bin.[find best split in the range of this bin];
    splitPriorityQueue.[add (split)];
    fringe.[delete (bin)];
  }
  nextBestSplit = splitPriorityQueue.[give best split in queue];
  newBinLeft, newBinRight =
    nextBestSplit.[perform split on its bin and replace the bin with the
    two new bins newBinLeft and newBinRight];
  numSplits++;
  fringe.[add the two new bins (newBinLeft, newBinRight)];
} UNTIL (numSplits == maxNumBins - 1);
    
```

Fig. 4. Pseudo code for the tree building algorithm

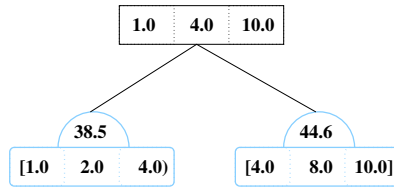


Fig. 5. Tree after the first cut

... a ... i a g b a d i i , b a d i i , h a i a c ... a i a i e e i e e i a e . We a ... h i e h d ... e i e d d i c e i a i ... i g b e ... d e e a i ... The ... d e d e f ... a g i h ... h ... i Fig e 4 .

I n t h e f ... i g e e e a e a e i g h e d a e f ... Fig e 1 a d 2 . F i l l t h e b e c ... i i f d i h e h e a g e a d ... e b i a e f ... e d . W i t h t h e ... b a g e ... e ... c a ... i a c ... i a e e a c h e d f ... B h ... i a e e a a e d a d a g - i e i h d f ... h e d i i i ... h e e ... i g h e e b i i c ... e d f ... b h ... i b e ... i . Fig e 5 . h ... h e d i c e i a i ... e e c ... e ... d i g ... h i i a i ... The ... d e e e e ... h e ... c a d h e ... e a f ... d e e e e ... h e e ... i b e c ...¹

E a c h ... d e e e e ... a b a g e , h e ... d e h e h e a g e . T h e a i - a b e ... i e ... h e e f a d i g h i d e f h e ... a e c ... e ... d i g ... a ... d e e e e ... h e ... i i ... a d ... a i ... f h e b a g e . T h e ... e a ... i i ... a d ... a i ... f h i e a ... e d a e a e 1.0 a d 10.0 . E a c h e a f ... d e e e -

¹ All values are rounded.

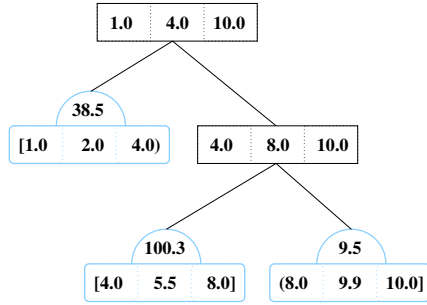


Fig. 6. Tree after the second cut

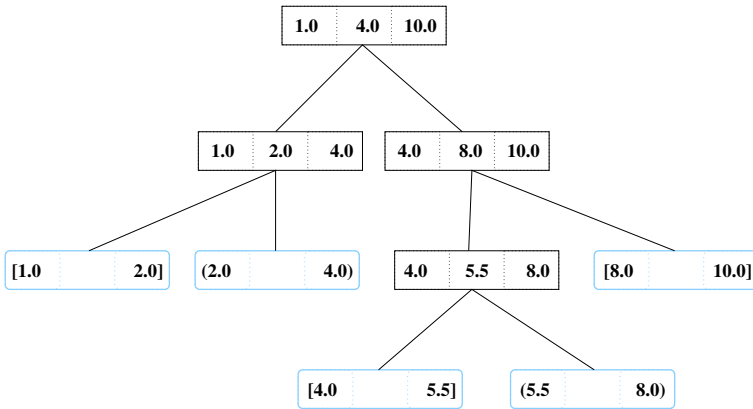


Fig. 7. Finalized tree

... a binary and hence a binary tree if the ... a ... of ...
 the binary data (if it is a ... of the ...). The ... a ... of ...
 a ... age. The ... a ... of ... the ... of ... the ... of ...

The ... the ... age ... the ... age 4.0. The ... the ...
 the ... 2.0 | 8.0. The ... the ... the ... data ... the ...
 [1.0:2.0] [2.0:4.0] [4.0:10.0] ... [1.0:4.0] [4.0:8.0] [8.0:10.0], ...
 each ... the ... the ... of ... the ... the ...
 ... The ... a ... 2.0, ... the ... age ... of 38.5 ...
 ... Fig. 6, the ... a ... 8.0 has ... age ... of 44.6. A ...
 ... the ... age ... age ... of ... the ...
 the ... a ... 8.0. Fig. 6. ... the ... of ...

After the ... 8.0 is ... the ... age ... and ...
 the ... the ... of ... The ... 5.5 and 9.9, ...
 ... of 100.3 and 9.5 ... For the ...
 ... the ... (... 2.0) and ... 5.5.

After ... the ... the ...
 The ... the ... Fig. 7. ...

density estimation. The algorithm generates the histogram at the beginning of the discretization. Figure 2.1 shows the histogram of each leaf node of the histogram. Each leaf node density is calculated.

5.3 The Stopping Criterion

The histogram is a function of the histogram. Based on the histogram, the algorithm generates the histogram of the data. The histogram is a function of the data. The histogram is a function of the data.

We use the 10-fold cross-validation method to determine the best histogram. We use the 10-fold cross-validation method to determine the best histogram. We use the 10-fold cross-validation method to determine the best histogram.

The algorithm is a function of the histogram. The algorithm is a function of the histogram. The algorithm is a function of the histogram.

5.4 A Problem: Small Cuts

The histogram is a function of the histogram. The histogram is a function of the histogram. The histogram is a function of the histogram.

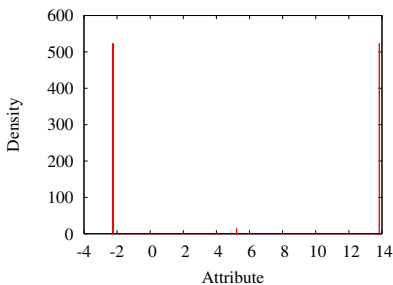


Fig. 8. Distorted histogram due to small cuts

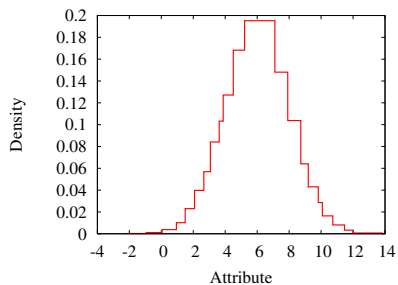


Fig. 9. Small cuts eliminated with heuristic

Here we are interested at the heuristic and the accuracy of the classification. More precisely, we are interested at the accuracy of the heuristic. The accuracy of the heuristic is defined as $\frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{h(x_i) = y_i\}}$, where N is the number of instances, h is the heuristic, and y_i is the true class label. The accuracy of the heuristic is high if the heuristic is able to correctly classify most of the instances. The accuracy of the heuristic is low if the heuristic is unable to correctly classify most of the instances.

6 Evaluation

We evaluated the TUBE decision algorithm on the 464 numeric attributes from 21 UCI datasets [7]. The algorithm was evaluated on the accuracy of the heuristic. The accuracy of the heuristic is defined as $\frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{h(x_i) = y_i\}}$, where N is the number of instances, h is the heuristic, and y_i is the true class label.

As a heuristic, we used the heuristic of the TUBE decision algorithm. The heuristic is defined as $h(x) = \text{argmax}_{c \in \mathcal{C}} \text{acc}_c(x)$, where $\text{acc}_c(x)$ is the accuracy of the heuristic for class c . The accuracy of the heuristic is defined as $\text{acc}_c(x) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{h(x_i) = c\}}$. The accuracy of the heuristic is high if the heuristic is able to correctly classify most of the instances. The accuracy of the heuristic is low if the heuristic is unable to correctly classify most of the instances.

Table 1. 464 numeric attributes from UCI datasets and their levels of uniqueness

Dataset	[0-20]	[20-40]	[40-60]	[60-80]	[80-100]	num inst
anneal	6	-	-	-	-	898
arrythmia	182	7	14	3	-	452
autos	13	-	-	1	1	205
balance-scale	4	-	-	-	-	625
winsconsin-breast-cancer	9	-	-	-	-	699
horse-colic	7	-	-	-	-	368
german-credit	6	-	-	-	1	1000
ecoli	7	-	-	-	-	336
glass	3	3	2	1	-	214
heart-statlog	12	1	-	-	-	270
hepatitis	4	1	1	-	-	155
hypothyroid	7	-	-	-	-	3772
ionosphere	2	-	2	31	-	351
iris	4	-	-	-	-	150
labor	8	-	-	-	-	57
lymphography	3	-	-	-	-	148
segment	14	3	-	2	-	2310
sick	7	-	-	-	-	3772
sonar	-	7	4	-	46	208
vehicle	17	1	-	5	-	846
vowel	-	-	4	8	-	990
Sum	315	23	27	51	48	
In percent	68	5	6	11	10	

The edge density of the resulting TUBE discretization is the edge density. The density of the edge is given by the edge weight 10^{-10} divided by the number of edges, which is 10^{-10} . The edge density is 10^{-10} .

On the other hand, the edge density (TUBE) is calculated against the edge density of the edge (EW-10), edge density of the edge (EWc B), edge density of the edge (EWc BO), and edge density of the edge (EF-10). The edge density of the edge (TUBE) is calculated against the edge density of the edge (EW-10), edge density of the edge (EWc B), edge density of the edge (EWc BO), and edge density of the edge (EF-10). The edge density of the edge (TUBE) is calculated against the edge density of the edge (EW-10), edge density of the edge (EWc B), edge density of the edge (EWc BO), and edge density of the edge (EF-10). The edge density of the edge (TUBE) is calculated against the edge density of the edge (EW-10), edge density of the edge (EWc B), edge density of the edge (EWc BO), and edge density of the edge (EF-10).

6.1 Evaluating the Fit to the True Distribution

Table 2 shows the results of the comparison. Each value in the table is the edge density of the edge (TUBE) against the edge density of the edge (EW-10), edge density of the edge (EWc B), edge density of the edge (EWc BO), and edge density of the edge (EF-10).

Table 2. Comparison of the density estimation results. Result of paired t-test based on cross-validated log-likelihood.

	EW-10	EWcvB	EWcvBO	EF-10
(0-20)				
TUBE significantly better	99	100	100	-
TUBE equal	1	0	0	-
TUBE significantly worse	0	0	0	-
[20-40)				
TUBE significantly better	48	43	43	48
TUBE equal	52	57	57	52
TUBE significantly worse	0	0	0	0
[40-60)				
TUBE significantly better	8	8	8	37
TUBE equal	92	92	92	63
TUBE significantly worse	0	0	0	0
[60-80)				
TUBE significantly better	53	56	56	67
TUBE equal	44	40	42	30
TUBE significantly worse	3	3	2	3
[80-100]				
TUBE significantly better	13	17	15	13
TUBE equal	85	81	81	85
TUBE significantly worse	2	2	4	2
Total				
TUBE significantly better	76	77	77	43
TUBE equal	23	22	22	55
TUBE significantly worse	1	1	1	2

ig 1 ca be e, e a . . . e e ec 1 e ba ed . . he c . . ec ed e a ed
 - e [8]. I a . . a ca e . . e h d 1 a ea a g da he he . e h d
 a d h . e ecia g d e . . 1 ca e 1 h . . 1 e e a d . . e ca e f
 high . 1 e e . We ha e a a ed he c . . e . . d i g a ib e a d he h
 ha TUBE 1 ge e a be e he a ib e e h ibi di c . 1 1 ie 1 he i
 di ib 1 . .

I 1 di c . . 1 he da a e . . ec 1 e 1 a ib e 1 h c . 1 . .
 di ib 1 . . a da ib e 1 h di c . 1 . . di ib 1 . . Da a e be . 20
 e ce . 1 e e ca be c . ide ed di c . 1 . . b he e a e . . e da a e .
 1 he hghe . 1 e e ca eg . ha h ed di c . 1 1 ie .

A ib e 1 h . . 1 e e e h ibi di c . 1 . . di ib 1 . . f di e -
 e . 1 d . S . e f he a ib e a e e di c e e a d ha e . . 1 ege a e
 (e.g. ehic e-9) . a . . ec 1 . (e.g. 1 1 -4) . . e ha e 1 eg a di ib ed
 da a 1 e (e.g. eg e -7) a d . e ha e da a 1 e 1 eg a 1 e a . (e.g.
 ba a ce- ca e-1). I he ca eg . f (0-20) 1 e e TUBE . . e f . . a
 . he . e h d . a . . a f he da a e .

I he ca eg . [60-80) ha f f he a ib e ha e a di ib 1 . ha 1 a
 1 . e be ee c . 1 . . da a a d di c e e da a (. . f he 1 . . he e a -

Table 3. Comparison of the number of bins

	EW-10	EWcvB	EWcvBO	EF-10
(0-20)				
TUBE significantly fewer	14	62	62	-
TUBE equal	2	8	7	-
TUBE significantly more	84	30	31	-
[20-40)				
TUBE significantly fewer	31	13	26	31
TUBE equal	4	30	17	4
TUBE significantly more	65	57	57	65
[40-60)				
TUBE significantly fewer	29	46	54	29
TUBE equal	38	42	38	38
TUBE significantly more	33	12	8	33
[60-80)				
TUBE significantly fewer	44	94	97	44
TUBE equal	14	6	3	14
TUBE significantly more	42	0	0	42
[80-100]				
TUBE significantly fewer	96	85	92	96
TUBE equal	2	15	8	2
TUBE significantly more	2	0	0	2
Total				
TUBE significantly fewer	29	65	68	56
TUBE equal	5	12	9	12
TUBE significantly more	66	23	23	32

lib e i hi ca eg, ha ea i ed di (ib i). TUBE' de i e i a i.
 a be e f, a he ea (ib e.

6.2 Comparing the Number of Bins

Table 3 shows a comparison of the performance of the different methods. The results show that the proposed method (TUBE) achieves a higher accuracy than the other methods. The results also show that the proposed method is more robust to noise than the other methods.

7 Conclusion

TUBE discretization provides a good approximation of the underlying data distribution. The results show that the proposed method (TUBE) achieves a higher accuracy than the other methods. The results also show that the proposed method is more robust to noise than the other methods.

Overall, the proposed method (TUBE) provides a good approximation of the underlying data distribution. The results show that the proposed method (TUBE) achieves a higher accuracy than the other methods. The results also show that the proposed method is more robust to noise than the other methods.

References

1. Eibe Frank and Ian H. Witten. Making better use of global discretization. *Proc. of the Sixteenth International Conference on Machine Learning*, pages 115–123, 1999.
2. B.W.Silverman. *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, 1986.
3. D.W.Scott. *Multivariate Density Estimation*. John Wiley & Sons, 1992.
4. Ian H. Witten and Eibe Frank. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann, 2000.
5. P. Smyth. Model selection for probabilistic clustering using cross-validated likelihood. *Statistics and Computing*, pages 63–72, 2000.
6. J.R.Quinlan. *C4.5 Programs for Machine Learning*. Morgan Kaufmann, 1993.
7. C.L. Blake S. Hettich and C.J. Merz. UCI repository of machine learning databases, 1998.
8. C. Nadeau and Y. Bengio. Inference for the generalization error. *Machine Learning*, 52:239–281, 2003.
9. Y. Yang and G. I. Webb. Proportional k-interval discretization for naive-bayes classifiers. In *Proceedings of the 12th European Conference on Machine Learning (ECML)*, pages 159–173, Tokyo, 2001.

Weighted Average Pointwise Mutual Information for Feature Selection in Text Categorization

Karl-Michael Schneider

Department of General Linguistics, University of Passau,
94030 Passau, Germany
schneide@phil.uni-passau.de

Abstract. Mutual information is a common feature score in feature selection for text categorization. Mutual information suffers from two theoretical problems: It assumes independent word variables, and longer documents are given higher weights in the estimation of the feature scores, which is in contrast to common evaluation measures that do not distinguish between long and short documents. We propose a variant of mutual information, called *Weighted Average Pointwise Mutual Information* (WAPMI) that avoids both problems. We provide theoretical as well as extensive empirical evidence in favor of WAPMI. Furthermore, we show that WAPMI has a nice property that other feature metrics lack, namely it allows to select the best feature set size automatically by maximizing an objective function, which can be done using a simple heuristic without resorting to costly methods like EM and model selection.

1 Introduction

Automatic text categorization, i.e. the assignment of text documents to predefined categories, is an important task in many NLP applications. The common *bag of words* approach results in a document space with very high dimensionality. In order to speed up parameter estimation and classification and to improve the classifier performance, it is common to use feature selection to reduce the dimensionality of the document space. This is typically done using a filtering approach [1] in which each feature is assigned a score based on an independent evaluation, and the features are then ranked according to their scores, and the N highest ranked features are selected, where N is the desired vocabulary size. Wrapper methods, which use the classifier directly to evaluate different feature subsets [1], are not commonly used for text classification because of the high dimensionality of the feature space that makes searching for the best feature subset intractable.

Mutual Information (MI) is an information-theoretic measure that is often used to evaluate features. It measures the amount of information that the value of a feature in a document (e.g. the presence or absence of a word) gives about the class of the document. Feature selection studies have obtained good results with MI [2]. However, there are two problems associated with the use of MI for feature ranking: First, MI treats each feature as an independent random variable. This is a problem because words in a text are not independent. Second, classifiers based on generative models, such as Naive

Bayes [3], estimate class-conditional probability distributions over words from training data. In the multinomial Naive Bayes model [3,4] this is done by concatenating the training documents in each class to one long document and estimating the distribution of words in this long document. This gives larger weights to longer documents. However, in classifier evaluation, all (test) documents have equal weight irrespective of their length—that is, there is a mismatch between classifier training and evaluation.

This paper proposes a variant of MI, *Weighted Average Pointwise Mutual Information* (WAPMI) that avoids both aforementioned problems. We present theoretical (using an information-theoretic argument that links WAPMI to multinomial Naive Bayes) and empirical evidence (through extensive experimentation) in favor of WAPMI. WAPMI improves the performance of multinomial Naive Bayes over MI on a variety of standard benchmark corpora. It also outperforms several other standard metrics for feature ranking.

In addition, WAPMI has a very nice property compared to other metrics, including MI: It allows to determine the (theoretically) best feature set size by maximizing an objective function. This can be done using a simple heuristic by applying a general, data-independent threshold to the feature scores, without the need to resort to computationally intensive methods like EM and model selection. Other feature metrics only evaluate the relative usefulness, and it is not entirely clear how they could be used to define an objective function for feature selection.

We demonstrate the effectiveness of this general thresholding method in our experiments. On some datasets (notably those that are commonly regarded “easy” classification tasks) we obtain smaller feature sets and better performance, while on “difficult” datasets (i.e. large datasets with great variability in the vocabulary) WAPMI selects larger feature sets than other metrics while outperforming them.

The paper is structured as follows. In Sect. 2 we review the probabilistic framework of multinomial Naive Bayes. In Sect. 3 we define weighted average pointwise mutual information and motivate its use for feature ranking. We also discuss its relation to distributional clustering. The experimental setup is described in Sect. 4, and Sect. 4 presents our experiments and the results. Section 5 finishes with some conclusions.

2 Naive Bayes

Naive Bayes is a simple probabilistic classifier that is widely used for text classification [3,4]. Despite this independence assumption, Naive Bayes performs surprisingly well on text classification problems [5].

Let $C = \{c_1, \dots, c_{|C|}\}$ denote the set of possible classes of documents, and let $V = \{w_1, \dots, w_{|V|}\}$ be a vocabulary. The multinomial Naive Bayes classifier assumes that a document d is drawn from a multinomial distribution by $|d|$ independent trials on a random variable $W \in V$ with class-conditional distribution $p(w_t | c_j)$ (where $|d|$ denotes document length):

$$p(d|c_j) = p(|d|)|d|! \prod_{t=1}^{|V|} \frac{p(w_t | c_j)^{x_t}}{x_t!}$$

x_t is the number of times W yields w_t , i.e. the number of times the word w_t occurs in d . The parameters $p(w_t|c_j)$ are usually estimated from training documents using maximum likelihood with Laplace smoothing to avoid zero probabilities:

$$\hat{p}(w_t|c_j) = \frac{1 + n(c_j, w_t)}{|V| + n(c_j)}$$

where $n(c_j, w_t)$ is the number of occurrences of w_t in the training documents in c_j and $n(c_j)$ is the total number of word occurrences in c_j .

The posterior probability of the class given the document is given by Bayes' rule:

$$p(c_j|d) = \frac{p(c_j)p(d|c_j)}{p(d)}$$

where $p(d)$ is the total probability of d :

$$p(d) = \sum_{j=1}^{|C|} p(c_j)p(d|c_j)$$

The class priors $p(c_j)$ are estimated from training documents as the fraction of documents in class c_j . Given a document, the Naive Bayes classifier selects the class with the highest posterior probability (we can omit those parts that do not depend on the class in the maximization):

$$c^*(d) = \operatorname{argmax}_{c_j} p(c_j)p(d|c_j) \quad (1)$$

3 Weighted Average Pointwise Mutual Information

3.1 Defining Weighted Average Pointwise Mutual Information

Mutual Information is a measure of the information that one random variable gives about the value of another random variable [6]. Let W be a random variable that ranges over the vocabulary V , and let C be random variable that ranges over classes. The mutual information between W and C is defined as:

$$I(W;C) = \sum_{t=1}^{|V|} \sum_{j=1}^{|C|} p(w_t, c_j) \log \frac{p(w_t|c_j)}{p(w_t)} \quad (2)$$

The term $\log \frac{p(w_t|c_j)}{p(w_t)}$ is called *pointwise mutual information* [7].¹ Note that mutual information can be written as a weighted sum of Kullback-Leibler (KL) divergences. The KL-divergence between two probability distributions p and q is defined as $D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$ [6]. Thus (2) can be written as the weighted average KL-divergence

¹ In [2] this is called *information gain*, and the term *mutual information* is used as a synonym for pointwise mutual information.

between the class-conditional distribution of words and the global (unconditioned) distribution in the entire corpus:

$$I(W;C) = \sum_{j=1}^{|C|} p(c_j) D(p(W|c_j) || p(W))$$

To rank features we would like a measure for each feature. A common method is to define new binary random variables, W_t , for each word that indicate whether the next word in a document is w_t (or some other word) [3,8]: $p(W_t = 1) = p(W = w_t)$. Then the MI-score for w_t is given by:

$$MI(w_t) := I(W_t;C) = \sum_{j=1}^{|C|} \sum_{x=0,1} p(W_t = x, c_j) \log \frac{p(W_t = x|c_j)}{p(W_t = x)} \quad (3)$$

The problem with (3) is that it treats W_t as an independent random variable, but in fact $\sum_{t=1}^{|V|} p(W_t = 1) = 1!$ To avoid this independence assumption, we consider (2) as a sum over word scores, where the score for w_t is the pointwise mutual information with the class, averaged over all classes:

$$PMI(w_t) := \sum_{j=1}^{|C|} p(w_t, c_j) \log \frac{p(w_t|c_j)}{p(w_t)} \quad (4)$$

The problem with (4) is that it treats all training documents in one class as one big document (because of the way the class-conditional probabilities are estimated). Thus, if there is variation in the document lengths, (4) is dominated by the longer documents. To avoid this problem, we replace the weight $p(w_t, c_j)$ with a term that is a weighted average of the document-conditional probabilities $p(w_t|d_i) = n(w_t, d_i)/|d_i|$ where $n(w_t, d_i)$ is the number of times w_t occurs in d_i and $|d_i|$ is the length of d_i .² Thus weighted average pointwise mutual information is defined as:

$$WAPMI(w_t) := \sum_{j=1}^{|C|} \sum_{d_i \in c_j} \alpha_i p(w_t|d_i) \log \frac{p(w_t|c_j)}{p(w_t)} \quad (5)$$

We consider several alternatives for the weights α_i , which can be associated with different measures for classifier evaluation:

- $\alpha_i = p(c_j) \cdot |d_i| / \sum_{d_i \in c_j} |d_i|$. This gives each document a weight proportional to its lengths and yields (4).
- $\alpha_i = 1 / \sum_{j=1}^{|C|} |c_j|$. This gives equal weight to all documents. This corresponds to an evaluation measure that counts each misclassified document as the same error, i.e. classification accuracy.
- $\alpha_i = 1 / (|c_j| \cdot |C|)$ where $d_i \in c_j$. This gives equal weight to the classes by normalizing for class size, i.e. documents from smaller categories receive higher weights. This compensates for the dominance of larger categories in classifier evaluation.

² Note that any word that does not occur in d_i has zero probability.

By summing (5) over all words we obtain the total weighted average pointwise mutual information between the word variable W and the class variable C :

$$WAPMI(W;C) := \sum_{t=1}^{|V|} \sum_{j=1}^{|C|} \sum_{d_i \in c_j} \alpha_i p(w_t | d_i) \log \frac{p(w_t | c_j)}{p(w_t)} \quad (6)$$

In the following subsections we provide theoretical evidence that total WAPMI could be used as an objective function, and the goal of feature selection is to maximize that objective function.

3.2 Relation to Distributional Clustering

Note that (6) can be written as a weighted sum of the difference between (i) the KL-divergence of the document-conditional distribution from the corpus distribution and (ii) the KL-divergence of the document-conditional distribution from the class-conditional distribution:

$$\sum_{j=1}^{|C|} \sum_{d_i \in c_j} \alpha_i \left[D(p(W|d_i) \| p(W)) - D(p(W|d_i) \| p(W|c_j)) \right] \quad (7)$$

This can be interpreted as an estimate of how similar the documents in one class are and how dissimilar documents of different classes are. From a clustering perspective we can say that (7) is large if the documents that belong to the same class form tight clusters, with wide separation between the clusters. Interpreting text categorization as an information retrieval task (i.e. regarding classes as queries) this is a desirable property that has been argued to improve document retrieval performance in the vector space model [9].

In distributional clustering the goal is to cluster similar objects (e.g. documents) together so as to maximize the value of an objective function that measures the quality of the clustering [10]. Below we argue that maximizing (7) is expected to improve the accuracy of the multinomial Naive Bayes classifier. Thus we can regard total weighted average pointwise mutual information as an objective function (since it is a function of the entire training corpus). However, in contrast to clustering, we do not change the clusters (which correspond to the classes in the training corpus and which we consider to be fixed). Instead our goal is to improve the clustering by changing the document representation (i.e. by using a subset of the features).

3.3 Relation to Multinomial Naive Bayes

We can use (7) to get an estimate of the expected performance of Naive Bayes on the training set (and by generalization also on a test set, if the test documents are drawn from the same distribution). We manipulate the Naive Bayes classifier (1) in an information theoretic framework using the fact that a document defines a probability distribution over words. We define the distance of a document, d_i , from a class, c_j , as the KL-divergence between the document-conditional word distribution and the class-conditional distribution. Naive Bayes can then be written in the following form by taking logarithms, dividing by the length of d_i and adding the entropy of d_i , $H(p(W|d_i)) = -\sum_t p(w_t | d_i) \log p(w_t | d_i)$ [10]:

$$c^*(d_i) = \arg \min_{c_j} \left[D(p(W|d_i) \| p(W|c_j)) - \frac{1}{|d_i|} \log p(c_j) \right] \quad (8)$$

Note that the modifications in (8) do not change the classification of documents. Assuming equal class priors, Naive Bayes can thus be interpreted as selecting the class which has the least distance from the document. Taking into account the arguments from the previous subsection, maximizing the total weighted average pointwise mutual information (6) would thus increase the probability that each document is nearer to its true class than to any other class, and would therefore be classified correctly by multinomial Naive Bayes.

3.4 Using WAPMI as an Objective Function for Feature Selection

Taking into account the arguments in the previous subsections, the best feature set would be one that maximizes the total WAPMI (6). Note that the WAPMI score (5) can be negative, which suggests the following simple heuristic for maximizing total WAPMI: Simply select all words with a positive WAPMI score and removing all other words. This is equivalent to applying a threshold of $\theta = 0$ to the WAPMI score. We examine this empirically in Sect. 4. In contrast, mutual information is always non-negative (and almost always positive), and it is not entirely clear how mutual information could be used as an objective function in feature selection.

Note that the above heuristic is only an approximation. In fact, feature selection isn't entirely well-defined in multinomial Naive Bayes, since we are not only pruning the model but the data too! Pruning the vocabulary changes the distribution of the remaining words. An alternative would be to not greedily discard words but perform several iterations and recompute the objective function after each iteration until convergence. We tried this, but there was almost no difference. In most cases, convergence occurred after only two or three iterations, with only a few additional words removed after the first round.

4 Experiments

4.1 Datasets and Procedures

We perform experiments on five text categorization datasets, described in Table 1. The 20 Newsgroups dataset³ consists of Usenet articles distributed evenly in 20 different newsgroups that make up the classes [11]. We remove newsgroup headers and binary attachments and use only words consisting of alphabetic characters as tokens, after converting to lower case and mapping numbers, URLs and email addresses to special tokens.

The WebKB dataset and the 7 Sectors dataset are both available from the WebKB project [12].⁴ WebKB contains web pages gathered from computer science departments and categorized in six classes plus one *other* class. We use only the four most populous classes *course*, *faculty*, *project* and *student*. The 7 Sectors data consists of web pages

³ <http://people.csail.mit.edu/people/jrennie/20Newsgroups/>

⁴ <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-11/www/wwwkb/>

Table 1. Corpus statistics. The last two columns show the number of documents in the smallest and biggest categories, respectively.

Dataset	Classes	Vocabulary	Documents	Smallest	Largest
20 Newsgroups	20	94,897	19,997	997	1,000
WebKB	4	41,015	4,199	504	1,641
7 Sectors	48	42,110	4,582	39	105
Reuters-10 (train)	10	22,430	6,490	181	2,877
Reuters-10 (test)	10	13,849	2,545	56	1,087
Reuters-90 (train)	90	24,719	7,770	1	2,877
Reuters-90 (test)	90	15,660	3,019	1	1,087

from different companies divided into a hierarchy of classes. We use the flattened version of the data. We strip all HTML tags and use only words and numbers as tokens, after converting to lower case and mapping numbers and other expressions to special tokens.

The Reuters-21578 dataset⁵ consists of Reuters news articles belonging to zero or more topic classes. We use the ModApte split [13] and produce two versions of the corpus. Reuters-10 uses only the 10 largest topics. On average, each document belongs to 1.105 topic classes. Reuters-90 uses all 90 topics that have at least one document in the training and test set, with an average of 1.235 topics per document.

Except on Reuters, all experiments are performed using cross-validation. We follow the methodology in [3]. For 20 Newsgroups and 7 Sectors, we split the data into five parts of equal size and with equal class distribution. For WebKB we produce ten train/test splits using stratified random sampling with 70% training and 30% test data. We report average classification accuracy across trials.

For the Reuters experiments we build a binary classifier for each topic, using the documents belonging to each topic as positive examples and all other documents as negative examples. Following the standard methodology with multi-label datasets, we ignore the classification decision of the classifier and use the classification scores to rank the documents. We then report precision/recall breakeven points averaged over all topics (called “macroaverage”). Instead of the Naive Bayes posterior probabilities, which tend to produce extreme values with growing document length due to the Naive Bayes independence assumption and are not comparable across documents, we use the normalized KL-divergence based classification scores described in [12].

4.2 Quality of Selected Features

We compare our WAPMI scoring function against three other scoring functions: Mutual Information [3], Chi-squared [2] and Bi-normal separation [14]. We evaluate the quality of the selected features by varying the number of selected features. We use WAPMI with equal weighting for all documents (we also experimented with equal class weights but found no statistically significant difference). Table 2 shows the top 20 words in the entire 20 Newsgroups corpus according to Mutual Information and WAPMI.

Figure 1 shows classification accuracy on the three datasets. As can be seen, the WAPMI scoring function yields higher classification accuracy, although on WebKB

⁵ <http://www.daviddlewis.com/resources/testcollections/reuters21578/>

Table 2. 20 words with highest MI (left) and WAPMI score (right) in the 20 Newsgroups corpus

MI	Word	MI	Word	WAPMI	Word	WAPMI	Word
0.02833	ax	0.00174	g	0.00221	rainbowthreedigit	0.00073	rainbowdigits
0.01555	rainbowonedigit	0.00168	w	0.00179	sale	0.00070	mac
0.00387	rainbowdigits	0.00161	m	0.00150	rainbowtwodigit	0.00068	clipper
0.00374	rainbowtwodigit	0.00155	u	0.00140	windows	0.00067	taggedemail
0.00336	x	0.00144	v	0.00129	x	0.00067	card
0.00222	q	0.00143	of	0.00091	car	0.00066	thanks
0.00188	rainbowthreedigit	0.00124	god	0.00089	god	0.00065	team
0.00182	f	0.00119	r	0.00087	game	0.00064	he
0.00181	max	0.00109	p	0.00083	drive	0.00064	i
0.00175	the	0.00104	that	0.00074	bike	0.00064	space

the difference is statistically significant only for up to 2,000 words. In general, the improvement seems to be higher on smaller vocabulary sizes.

The class distribution is highly skewed in the Reuters datasets. The largest category (earn) has 2,877 documents in the training set, while the smallest category in Reuters-10 (corn) has 181 documents in the training set. In Reuters-90 there are 29 categories with less than 10 documents in the training set.

For the Reuters experiments we use two versions of WAPMI: with equal weights for all documents (WAPMI1), and with equal class weights (WAPMI2) (cf. Sect. 3.1), which deemphasizes the impact of the larger classes. Figure 2 shows the results on the Reuters datasets with 10 and 90 categories. We report macroaveraged precision/recall breakeven, which gives equal weight to the performance on each category. WAPMI with equal weights on documents does not perform better than the other metrics, except for very small vocabularies on Reuters-90. However, when the weights are set such that documents from smaller categories receive higher weights (WAPMI2), WAPMI clearly outperforms the other feature scoring methods.

4.3 Global Thresholding

In addition to the experiments with varying numbers of features we also examined the possibility of using a global thresholding strategy, with a fixed threshold that is applied to all datasets. We are interested in the sensitivity of the various feature scoring functions to the difficulty of the classification task. In general, the Naive Bayes classifier performs better with large vocabularies, but the optimal vocabulary size depends on the dataset. For instance, the 20 Newsgroups dataset requires a larger vocabulary for optimal classification accuracy than the other datasets [3].

For Mutual Information, Chi-squared and Bi-normal separation we select a threshold that yields relatively good performance on all datasets. For WAPMI we use the theoretically best threshold 0. For all datasets except 20 Newsgroups we use both variants with equal weights on documents (WAPMI1) and on classes (WAPMI2). For 20 Newsgroups WAPMI1 and WAPMI2 are the same because all classes have the same number of documents.

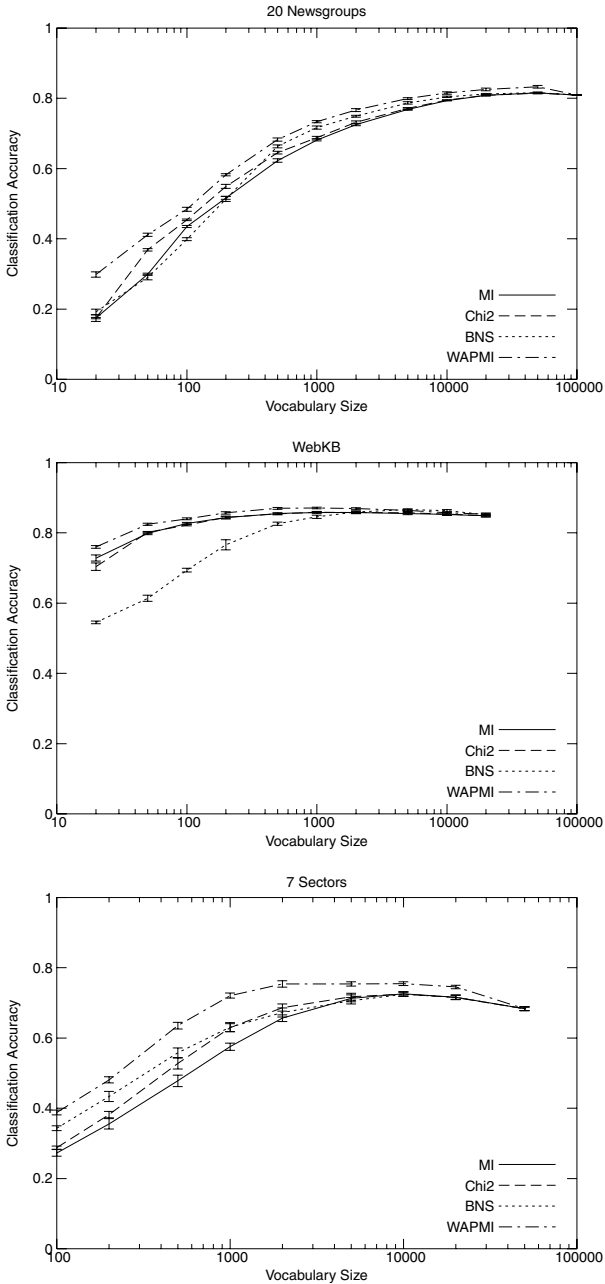


Fig. 1. Classification accuracy on 20 Newsgroups (top), WebKB (middle) and 7 Sectors (bottom). Curves show small error bars twice the width of the standard error of the mean. Differences between WAPMI and the other metrics are statistically significant (at the 95% confidence level using a two-tailed paired t-test) at the following vocabulary sizes: on 20 Newsgroups from 20 to 50,000 words; on WebKB from 20 to 2,000 words; on 7 Sectors from 100 to 20,000 words.

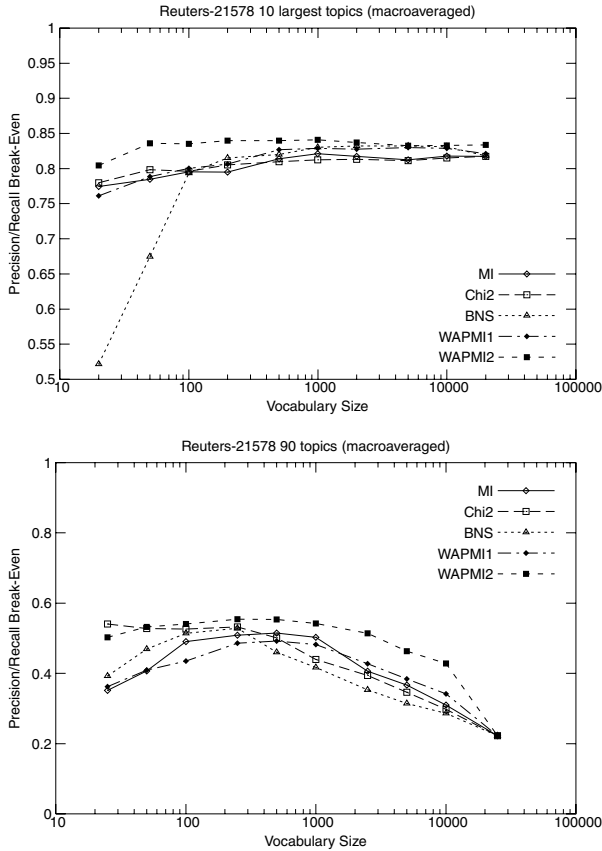


Fig. 2. Macroaveraged precision/recall breakeven on the Reuters datasets with 10 (top) and 90 (bottom) topic classes. WAPMI1 gives equal weight to documents, while WAPMI2 gives equal weight to classes.

Table 3 shows the results. For each dataset and each scoring function we report the number of features and the classification performance at the selected threshold. In addition we show the classification performance at the full vocabulary (i.e. with no feature selection).

We make two observations in Table 3. First, WAPMI is always among the top performers, although its performance is significantly better only on 20 Newsgroups and Reuters. Mutual Information performs significantly worse than the other metrics on 7 Sectors. Secondly and more importantly, the number of features selected by WAPMI seems to reflect the difficulty of the datasets better than for the other scoring methods. For 20 Newsgroups, which requires many features, WAPMI1 selects more features than any other method, while it still omits some features which results in an improvement of 2 percentage points compared to the full vocabulary. In contrast, the WAPMI scores select considerably less features on the Reuters datasets than the other methods, with better results.

Table 3. Global thresholding results. Shown are the number of selected words at the predefined threshold, classification performance, and standard deviation where applicable. Statistically significant differences (at $p = 0.95$ using a two-tailed paired t-test) are printed in boldface. For Reuters, macroaveraged precision/recall breakeven points are shown.

	20 Newsgroups			WebKB			7 Sectors		
	Words	Acc	SDev	Words	Acc	SDev	Words	Acc	SDev
$\text{Chi}^2=0.1$	65,194	81.35%	0.36%	32,712	84.79%	0.99%	15,147	72.29%	1.19%
$\text{MI}=10^{-7}$	77,694	81.13%	0.37%	32,776	84.79%	1.01%	37,474	68.32%	1.17%
$\text{BNS}=0.05$	62,777	81.42%	0.26%	32,550	84.78%	0.99%	8,545	72.27%	1.78%
$\text{WAPMI1}=0$	85,870	82.92%	0.72%	32,091	85.00%	0.96%	37,422	73.12%	0.57%
$\text{WAPMI2}=0$				32,278	85.06%	1.03%	37,428	73.14%	1.01%
Full	86,019	80.97%	0.29%	32,873	84.80%	0.99%	37,474	68.32%	1.17%

	Reuters-10		Reuters-90	
	Words	P/R	Words	P/R
$\text{Chi}^2=0.1$	18,861	81.72%	23,395	22.30%
$\text{MI}=10^{-7}$	18,014	81.72%	22,571	22.57%
$\text{BNS}=0.05$	20,086	81.76%	23,778	22.26%
$\text{WAPMI1}=0$	7,617	82.47%	3,066	44.58%
$\text{WAPMI2}=0$	10,610	83.17%	20,762	38.97%
Full	22,430	81.61%	24,719	22.28%

5 Conclusions

This paper proposes weighted average pointwise mutual information (WAPMI) as a replacement for mutual information to rank features for feature selection in text categorization. Experiments on a number of standard benchmark datasets show that WAPMI outperforms several other feature scoring metrics, including mutual information, Chi-squared and Bi-normal separation. An important property of WAPMI is that the feature set size (i.e. the number of selected features) can be set automatically, depending on the complexity and difficulty of the dataset, by using a simple constant-threshold heuristics that maximizes an objective function and does not require EM or model selection.

WAPMI contains weights that can be set to account for skewed class distributions, which we used in our experiments with the Reuters dataset and obtained improved classification performance. It is not entirely clear how this could be done with other metrics.

We have used WAPMI with the multinomial Naive Bayes classifier, but future work should deal with other classification models, e.g. support vector machines. A general open problem is that feature selection for multinomial Naive Bayes is not entirely well-defined, thus we are actually approximating feature selection. More work is required to better understand how feature selection affects the class-conditional distributions.

Acknowledgments

The author would like to thank the anonymous reviewers for their detailed comments and suggestions that helped to improve the paper.

References

1. John, G.H., Kohavi, R., Pfleger, K.: Irrelevant features and the subset selection problem. In Cohen, W.W., Hirsh, H., eds.: *Machine Learning: Proceedings of the Eleventh International Conference*, San Francisco, CA, Morgan Kaufmann Publishers (1994) 121–129
2. Yang, Y., Pedersen, J.O.: A comparative study on feature selection in text categorization. In: *Proc. 14th International Conference on Machine Learning (ICML-97)*. (1997) 412–420
3. McCallum, A., Nigam, K.: A comparison of event models for Naive Bayes text classification. In: *Learning for Text Categorization: Papers from the AAAI Workshop*, AAAI Press (1998) 41–48 Technical Report WS-98-05.
4. Eyheramendy, S., Lewis, D.D., Madigan, D.: On the Naive Bayes model for text categorization. In Bishop, C.M., Frey, B.J., eds.: *AI & Statistics 2003: Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*. (2003) 332–339
5. Friedman, J.H.: On bias, variance, 0/1-loss, and the curse-of-dimensionality. *Data Mining and Knowledge Discovery* **1** (1997) 55–77
6. Cover, T.M., Thomas, J.A.: *Elements of Information Theory*. John Wiley, New York (1991)
7. Church, K.W., Hanks, P.: Word association norms, mutual information, and lexicography. *Computational Linguistics* **16** (1990) 22–29
8. Rennie, J.D.M.: *Improving multi-class text classification with Naive Bayes*. Master's thesis, Massachusetts Institute of Technology (2001)
9. Salton, G., Wong, A., Yang, C.S.: A vector space model for automatic indexing. *Communications of the ACM* **18** (1975) 613–620
10. Dhillon, I.S., Mallela, S., Kumar, R.: A divisive information-theoretic feature clustering algorithm for text classification. *Journal of Machine Learning Research* **3** (2003) 1265–1287
11. Lang, K.: NewsWeeder: Learning to filter netnews. In: *Proc. 12th International Conference on Machine Learning (ICML-95)*, Morgan Kaufmann (1995) 331–339
12. Craven, M., DiPasquo, D., Freitag, D., McCallum, A., Mitchell, T., Nigam, K., Slattery, S.: Learning to construct knowledge bases from the World Wide Web. *Artificial Intelligence* **118** (2000) 69–113
13. Apté, C., Damerau, F., Weiss, S.M.: Towards language independent automated learning of text categorization models. In: *Proc. 17th ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '94)*. (1994) 23–30
14. Forman, G.: An extensive empirical study of feature selection metrics for text classification. *Journal of Machine Learning Research* **3** (2003) 1289–1305

Non-stationary Environment Compensation Using Sequential EM Algorithm for Robust Speech Recognition *

Haifeng Shen, Jun Guo¹, Gang Liu¹, and Qunxia Li²

¹ Beijing University of Posts and Telecommunications, 100876, Beijing, China
shen_hai_feng@126.com, guojun@bupt.edu.cn, lg@pris.edu.cn

² University of Science and Technology Beijing, 100083, Beijing, China
kellylqx@163.com

Abstract. The paper presents a non-stationary environment compensation using sequential EM estimation for tracking the complicated environment. All of the noisy features used in the recognition system are effectively compensated. The speech corruption in the log domain such as the 24 log-filterbank coefficients and the log-energy feature can be modeled as a nonlinear model. For efficient estimating noise parameter using the subsequent sequential Expectation-Maximization (EM) algorithm, the nonlinear environment model is linearized by the truncated first-order vector Taylor series (VTS) approximation. Due to the cepstral features are nearly independence, we train the clean speech using cepstral features and the log-energy feature, and then obtain a diagonal Gaussian mixture model in the log domain by taking inverse discrete cosine transform (IDCT). The experiments are conducted on the large vocabulary continuous speech recognition (LVCSR) system. Results demonstrate that it achieves attractive improvements when compared with CMN (cepstral mean normalization) and the batch-EM based compensation approach.

1 Introduction

The recognition performance will be severely degraded in the acoustic-distorted environments due to mismatches between the training and the test environments. The test utterances represent specific conditions such as specific speakers, specific speaking styles, specific noisy conditions, which generally are not included in the training data set and usually differ from the training conditions. There are many compensation approaches for reducing the influences of these mismatches on the speech. CMN (cepstral mean normalization), with the merits of inexpensive computation load and good recognition performance, can remove the cepstral mean from all vectors with the cepstral mean calculated separately from each sentence assuming that the average cepstral mean in the training and testing environments are equal to each other. The data-driven approach such that SNR-Dependent Cepstral Normalization (SDCN), Fixed Codeword-Dependent Cepstral Normalization (FCDCN) [1], needs a “stereo”

* This research was sponsored by NSFC (National Natural Science Foundation of China) under Grant No.60475007, the Foundation of China Education Ministry for Century Spanning Talent and BUPT Education Foundation.

database that contains time-aligned samples of speech which had been simultaneously recorded in both the training and the comprehensive test environments. The cepstral features of the incoming speech are compensated by direct comparison. The problem of the data-driven approach is that the stereo data recorded in a specific test environment is not suitable for another real environment. Moreover, this kind of the approaches is really complicated in recording the “stereo” databases and not effective when dealing with the non-stationary environment. Recently, the model-based approach becomes the most attractive technique [2]-[10]. The acoustic-distorted environment is modeled as an explicit model. For effectively modeling the statistical distribution of the noisy observation and estimating the environment parameters, the environment model is postprocessed to achieve the compact model. For instance, by employing the truncated first order vector Taylor series (VTS) approximation [2] [3] [5] and statistical linear approximation (SLA) [6]-[8], the nonlinear model is linearized. It is proven that such environment approximation approaches achieve the considerable performance on speech recognition. Furthermore, based on maximum likelihood estimation (ML) [7] or maximum a posteriori estimation (MAP) criterion [3], the noise parameter is iteratively updated to the real value using EM algorithm, generally, using the batch-EM algorithm. It is clear that the batch-EM algorithm can be carried out assuming that the environment is stationary, that is, the noise statistics is iteratively updated by computing the posteriori probabilities of all of the incoming speech frames. Although the batch-EM algorithm also improves the recognition performance in the non-stationary environment, this improvement is rather limited, especially in the high time-varying environment. The sequential EM algorithm [8]-[10], can deal with this problem and can improve recognition performance considerably compared with the batch EM environment compensation, especially in the time-varying environment.

In this paper, we present a non-stationary environment compensation based on sequential EM algorithm. Generally speaking, most of state-of-the-art speech recognition systems use the Mel frequency cepstral coefficients (MFCCs) and the log-energy feature as the acoustic vector. It is well known that the log-energy feature also makes a significant contribution for improving the recognition performance. Because the cepstral coefficients can be obtained from the log-filterbank coefficients by taking DCT transform, a number of the papers in literature [2]-[6], [8] [10] deal with the cepstral coefficients or the log-filterbank coefficients for making the feature robust against the noise environments. But it is clear that if the log-energy feature isn't well compensated, the system also can deteriorate the system performance, especially in the condition with a large amount of noise. Therefore, in this paper, we compensate all of the log-filterbank coefficients and the log-energy feature. Then taking DCT transform and corresponding dynamic features computation, the compensated cepstral coefficients and the log-energy feature plus the first and second differentials are obtained. For effectively estimating the environment parameter, the environment in the log domain can be modeled as a nonlinear model and linearized using the truncated first-order VTS approximation. It is noticeable that the clean speech model has a severe influence on the recognition performance. Generally, the clean speech model is modeled as the diagonal Gaussian mixture distribution for the effectiveness of the subsequent environment parameter estimation, also for decreasing the huge computation

load. Due to the aforementioned nearly independence in the cepstral domain, our approach to this is based on combination of all cepstral coefficients and the log-energy feature. Then the diagonal clean model in the log domain is obtained by taking inverse DCT transform on the cepstral statistics of the trained model. Based on initializing the truncated first-order VTS coefficients by employing the current estimated noise parameter and the next noisy speech frame, we update the next frame noise parameter by using sequential EM algorithm until the last noisy frame. The experiments are conducted on the large vocabulary continuous speech recognition (LVCSR) system. Results demonstrate that the environment compensation using the sequential EM algorithm improves recognition performance considerably compared with the batch EM environment compensation, especially in the time-varying environment. After introducing the forgetting factor for tracking the non-stationary time-varying environment, the performance of the speech recognition can further be improved in the non-stationary environment. The rest of the paper is organized as follows. The next section briefly describes the environment model approximation and accordingly investigates the statistical characteristics of the noisy speech. In section 3, we present sequential EM algorithm for noise parameter estimation. The experimental results are given in section 4 and some conclusions are drawn in section 5.

2 Environment Model Approximation

As seen in the appendix, due to the noise is additive in the linear spectral domain, the speech corruption will be nonlinear in the log spectral domain. In addition, the log-energy feature has the same corruption form as those of the log filterbank coefficients. So we can describe the corruption of these features in the noisy environment jointly. Denote the noisy feature, the clean feature and the noise in the log domain by y , x and n . The corruption is well represented as

$$y = x + \log(1 + \exp(n - x)) = x + f(x, n). \quad (1)$$

We assume the clean speech is modeled as a Gaussian mixture model:

$$p(x) = \sum_{j=1}^M p_j N(x; \mu_{xj}, \Sigma_{xj}), \quad (2)$$

in which M denotes the number of mixture components, p_j , μ_{xj} and Σ_{xj} denote the mixture coefficient, the mean vector and the diagonal covariance matrix for the j th mixture component, respectively. In our system, we first train the clean cepstral coefficients and the log-energy feature to obtain Gaussian mixture model. Then taking inverse DCT transform on these cepstral probability statistics, Gaussian mixture model in the log domain can be derived. We assume the noise is a Gaussian and statistically independent from the clean speech. The probability distribution of the noisy speech, unfortunately, is not the Gaussian mixture model due to the nonlinear relationship between the noisy speech and the clean speech described in Eq.(1). To simplify the distribution of the noisy speech and efficient noise estimation using sequential EM algorithm in the following step, we employ the truncated first-order VTS

expansion to linearize the nonlinearity $f(x, n)$ in Eq.(1) around the vector points (μ_{xj}, n_0) . This gives the linearized model in the j th mixture component:

$$y = A_j x + B_j n + C_j, \tag{3}$$

where

$$\begin{cases} A_j = 1 + \nabla_x f(\mu_{xj}, n_0) \\ B_j = \nabla_n f(\mu_{xj}, n_0) \\ C_j = f(\mu_{xj}, n_0) - \nabla_x f(\mu_{xj}, n_0)\mu_{xj} - \nabla_n f(\mu_{xj}, n_0)n_0 \end{cases}, \tag{4}$$

and the gradients $\nabla_x f(\mu_{xj}, n_0)$ and $\nabla_n f(\mu_{xj}, n_0)$ have the following close form:

$$\begin{cases} \nabla_x f(\mu_{xj}, n_0) = \text{diag} \left(\frac{1}{1 + \exp\{n_0 - \mu_{xj}\}} \right) \\ \nabla_n f(\mu_{xj}, n_0) = 1 - \nabla_x f(\mu_{xj}, n_0) \end{cases} \tag{5}$$

3 Noise Estimation Using Sequential EM Algorithm

Assuming that the noise is a single Gaussian distribution with mean vector n_t and covariance matrix Σ_n in each instant time t , we can see that the distribution of the noisy speech is a Gaussian mixture model by applying the first-order VTS approximation. In this paper, for simplicity, we are only interested in the noise mean estimation in each frame. The covariance matrix of each frame is set with equal value and can be estimated from silence frames. Given the acoustic-distorted feature sequence $Y_{t+1} = \{y_1, y_2, \dots, y_{t+1}\}$ and the previous noise estimate sequence $\Lambda_m = \{\hat{n}_0, \hat{n}_1, \dots, \hat{n}_t\}$ in which \hat{n}_0 is the initial parameter estimate and \hat{n}_t is the noise estimate at time t , the noise \hat{n}_{t+1} at time $t+1$ can be obtained under ML criterion:

$$\hat{n}_{t+1} = \arg \max_{n_{t+1}} \{ \log P(Y_{t+1}, J_{t+1} | n_{t+1}, \Lambda_m) \}, \tag{6}$$

where $J_{t+1} = \{j_1, j_2, \dots, j_{t+1}\}$ is the a set of the mixture components up to time $t+1$.

In general, it is not easy to estimate instant noise parameter. In this section, we use the sequential EM algorithm to iteratively estimate the different instant noise. At each iteration, the likelihood in Eq.(6) are increase until convergence. The auxiliary function is given below

$$Q(\hat{n}_{t+1} | n_{t+1}, \Lambda_m) = E \{ \log P(Y_{t+1}, J_{t+1} | \hat{n}_{t+1}, \Lambda_m) | Y_{t+1}, n_{t+1}, \Lambda_m \}, \tag{7}$$

where n_{t+1} is the initial value needed to know beforehand. In the slow time-varying acoustic-distorted environment, the value n_{t+1} can be approximated using the previous estimate \hat{n}_t , then the above equation can be compactly written as

$$\begin{aligned}
 Q(\hat{n}_{t+1} | \Lambda_m) &\approx E\{\log P(Y_{t+1}, J_{t+1} | \hat{n}_{t+1}, \Lambda_m) | Y_{t+1}, \Lambda_m\} \\
 &\propto -\sum_{\tau=1}^{t+1} \sum_{j=1}^M p(j_\tau = j | y_\tau, \hat{n}_{\tau-1}) \{y_\tau - \hat{\mu}_{y_\tau, j}(\hat{n}_\tau)\}' \Sigma_{y_\tau, j}^{-1} \{y_\tau - \hat{\mu}_{y_\tau, j}(\hat{n}_\tau)\},
 \end{aligned} \tag{8}$$

where

$$\begin{cases} \hat{\mu}_{y_\tau, j}(\hat{n}_\tau) = A_j(\hat{n}_{\tau-1})\mu_{sj} + B_j(\hat{n}_{\tau-1})\hat{n}_\tau + C_j(\hat{n}_{\tau-1}) \\ \Sigma_{y_\tau, j} = A_j(\hat{n}_{\tau-1})\Sigma_{sj}A_j'(\hat{n}_{\tau-1}) + B_j(\hat{n}_{\tau-1})\Sigma_n B_j'(\hat{n}_{\tau-1}) \end{cases} \tag{9}$$

where the coefficients $A_j(\cdot)$, $B_j(\cdot)$ and $C_j(\cdot)$ are the functions of the noise paramter $\hat{n}_{\tau-1}$. That is, the nonlinear function $f(x, \hat{n}_\tau)$ in Eq.(1) is approximated around the vector point $(\mu_{sj}, \hat{n}_{\tau-1})$ by using vector Taylor expansion.

The posteriori probability $p(j_\tau = j | y_\tau, \hat{n}_{\tau-1})$ in Eq.(8) can be computed as

$$p(j_\tau = j | y_\tau, \hat{n}_{\tau-1}) = \frac{p_j N(y_\tau; \hat{\mu}_{y_\tau, j}, \Sigma_{y_\tau, j})}{\sum_{j=1}^M p_j N(y_\tau; \hat{\mu}_{y_\tau, j}, \Sigma_{y_\tau, j})} \tag{10}$$

where $\hat{\mu}_{y_\tau, j} = \mu_{sj} + f(\mu_{sj}, \hat{n}_{\tau-1})$.

In the non-stationary environment, the history observation data is not useful or not really important to current noise estimation. We can add the different weights according to their contributions on current noise estimation. The different weights can be added by introducing the forgetting factor ρ where ρ is a non-negative constant with value less than 1, thus, Eq.(8) can be rewritten as

$$Q(\hat{n}_{t+1} | \Lambda_m) = -\sum_{\tau=1}^{t+1} \rho^{t+1-\tau} \cdot \left\{ \sum_{j=1}^M p(j_\tau = j | y_\tau, \hat{n}_{\tau-1}) \{y_\tau - \hat{\mu}_{y_\tau, j}(\hat{n}_\tau)\}' \Sigma_{y_\tau, j}^{-1} \{y_\tau - \hat{\mu}_{y_\tau, j}(\hat{n}_\tau)\} \right\}. \tag{11}$$

By Taylor series expansion to the above auxiliary function, choosing the truncated second order items, and maximizing the approximated items with respect to the noise parameter, the noise at time $t + 1$ can be estimated [8]-[11]

$$\hat{n}_{t+1} = \hat{n}_t + \gamma \cdot \{K_{t+1}(\hat{n}_t)\}^{-1} S_{t+1}(\hat{n}_t), \tag{12}$$

where the disturbing factor γ is a non-negative constant with value greater than 0, the Fisher information matrix $K_{t+1}(\hat{n}_t)$ and the score vector $S_{t+1}(\hat{n}_t)$ are defined as following

$$\begin{aligned}
 K_{t+1}(\hat{n}_t) &= -\left. \frac{\partial^2 Q(n | \Lambda_m)}{\partial^2 n} \right|_{n=\hat{n}_t} = \sum_{\tau=1}^{t+1} \rho^{t+1-\tau} \sum_{j=1}^M p(j_\tau = j | y_\tau, \hat{n}_{\tau-1}) \cdot B_j'(\hat{n}_{\tau-1}) \Sigma_{y_\tau, j}^{-1} B_j(\hat{n}_{\tau-1}) \\
 &= \rho \cdot K_t + \sum_{j=1}^M p(j_{t+1} = j | y_{t+1}, \hat{n}_t) B_j'(\hat{n}_t) \Sigma_{y_{t+1}, j}^{-1} B_j(\hat{n}_t),
 \end{aligned} \tag{13}$$

$$S_{t+1}(\hat{n}_t) = \left. \frac{\partial Q(n | \Lambda_m)}{\partial n} \right|_{n=\hat{n}_t} = \sum_{j=1}^M p(j_{t+1} = j | y_{t+1}, \hat{n}_t) B_j(\hat{n}_t)' \Sigma_{y_{t+1}, j}^{-1} \{y_t - \hat{\mu}_{y_t, j}(\hat{n}_t)\}. \tag{14}$$

4 Experimental Results

A continuous hidden Markov model (HMM)-based speech recognition system is used in the recognition experiments for examining the presented approach. The utterances of 82 speakers (41 males and 41 females) from the mandarin Chinese corpus provided by the 863 plan (China High-Tech Development Plan[12]) are trained for triphone-based HMM acoustic models, where each triphone unit was modeled as a three-emitting-state left-right topology with a mixture of 16 Gaussian per state and diagonal covariance matrices. The utterances of 9 speakers from the clean corpus are used for subsequent artificial contamination with different noise class.

In order to extract Mel frequency cepstral coefficients (MFCCs) from the 16Hz noisy speech data, we use a power spectrum which is calculated every 10ms on a 25ms Hanning window with pre-emphasis coefficient 0.97, then take a mel-scaled triangular filterbank and logarithmic computation and accordingly obtain the Mel-scaled 24 log-filterbank coefficients. After transforming them into the cepstral domain with DCT transform, we obtain the first 12 cepstral coefficients (excluding the zero coefficients). The log-energy feature in each frame is computed after taking Hanning windowing. Accordingly, 39 dimensional features consisting of the 12 cepstral coefficients, the log-energy feature coefficient and their time derivatives are computed.

In our feature compensation paradigm, for modeling the clean speech, we extract a set of 24 MFCCs and one log-energy feature from the clean speech data for training and obtain a mixture of 128 Gaussian distributions. Then the mean vector of each mixture component in the Mel-scaled log spectral domain is obtained using inverse cosine transformation matrix. The covariance matrix is computed also from the cepstral domain using the inverse cosine transformation matrix and its transpose. By ignoring the off-diagonal elements in the covariance matrices assuming that the different coefficients are statistically independent, we obtain the diagonal covariance matrices in the log domain. With the developed sequential EM algorithm, the 24 dimensional log-filterbank features and a log-energy feature are compensated. With DCT transform and delta and delta-delta regression equations, the static coefficients (12 MFCCs plus the log-energy feature) and the corresponding dynamic coefficients (13 delta coefficients and 13 delta-delta coefficients) are computed.

In order to test the validity of the feature compensation algorithm, a number of experiments have been performed. They include the baseline without compensation, compensation with CMN (cepstral mean normalization), batch-EM estimation and sequential EM estimation. The forenamed three approaches are titled as “baseline”, “CMN” and “batch-EM”, respectively in Table 1 and Table 2. In the sequential EM estimation, to investigate the behavior in the non-stationary environment, we get three forms: “Seq-0.90”, “Seq-0.95” and “Seq-1.00” according to the different forgetting value ρ with 0.90, 0.95 and 1.00. And we add the stationary white noise and the non-stationary babble noise from NoiseX92 [13] to the test set according to different SNR varying from 0dB to 20dB. It is observed from Table 1 that, the sequential estimation gives considerable performances, compared with “baseline”, “CMN” and “Batch-EM”. For example, in the 5dB white noisy condition, “baseline” only achieves 2.54% recognition rate, “CMN” achieves 10.98% recognition rate, and “Batch-EM” achieves 17.00% recognition rate. The sequential estimation with the forgetting factor ρ set to 0.90, 0.95 and 1.00 gives 18.93%, 18.97% and 18.91% recognition rates and achieves

1.93%, 1.97%, and 1.91% improvements over that by “Batch-EM”, respectively. As a whole, the sequential estimation with different forgetting factor value achieves 0.77%, 0.75%, and 0.82% improvements over that by “Batch-EM”, respectively. It is clear that the presented approach is very effective in the stationary noisy condition.

Table 1. Recognition rates in the white noisy environment (%)

SNR	0dB	5dB	10dB	15dB	20dB	Avg.
baseline	0.32	2.54	11.14	30.00	56.04	20.01
CMN	3.31	10.98	29.99	36.94	61.65	28.57
Batch-EM	5.51	17.00	39.37	62.36	77.12	40.27
Seq-0.90	4.99	18.93	39.99	63.19	78.10	41.04
Seq-0.95	4.99	18.97	40.19	62.85	78.12	41.02
Seq-1.00	5.02	18.91	40.38	63.02	78.14	41.09

Table 2. Recognition rates in the babble noisy environment (%)

SNR	0dB	5dB	10dB	15dB	20dB	Avg.
baseline	3.87	24.61	54.84	62.98	80.23	45.31
CMN	11.16	32.38	55.95	71.04	80.31	50.17
Batch-EM	15.86	39.25	61.78	75.23	81.09	54.64
Seq-0.90	17.72	40.71	62.56	75.53	81.07	55.52
Seq-0.95	17.50	40.70	62.49	75.50	81.21	55.48
Seq-1.00	17.44	40.31	62.60	75.51	81.46	55.46

To test the validity of the sequential estimation in non-stationary conditions, we further test the babble noise in different SNR levels. It is observed in Table 2 that, using “baseline”, performance degradation is not obvious in the high SNR condition, such as in the 20dB condition. But when the noise amount increases, recognition performance quickly deteriorates with only 3.87% recognition rate in the 0dB condition. With “CMN” and “Batch-EM”, the phenomena can be relatively restrained. However, they still have the main limitations to cope with the non-stationary environments. Although compensation is applied to reduce the mismatch among the clean acoustic model and the test set, they remain a minor mismatch which they don’t obtain the best performance at all of non-stationary noisy conditions. With the sequential estimation algorithm, it can further reduce the mismatch and can improve the system performance in most of noisy conditions, especially in low SNR conditions. For example, in 5dB condition, the sequential EM algorithm with different forgetting factor achieves 1.46%, 1.45% and 1.06% improvements in comparison to “Batch-EM”, respectively.

From Table 1 and Table 2, we also observe that the sequential estimation averagely provides slight improvement when the forgetting factor ρ is 1.00 over that of ρ is 0.9 or 0.95 for the white noise. But we notice that it averagely provides slight improvement when ρ is 0.90 over that of ρ is 0.95 or 1.00 for the babble noise. The cause of this behavior is that the white noise is the stationary noise and the babble

noise is the non-stationary noise. For the white noisy condition, it is clear that the history data is very useful to noise estimation. With the reasonable forgetting factor, the presented approach can ignore the history data which is effective for computing the current noise parameter in the non-stationary condition. Due to the babble noise is a slow time-varying noise, the forgetting factor can be set with a high value relatively. For the highly time-varying conditions, ρ can be a low value to reasonably track the non-stationary characteristics.

5 Conclusions

We have presented an approach to environment compensation for robust speech recognition based on a sequential EM algorithm. The algorithm compensates entirely all of the features to deal with the environment corruption. The corruption causing distortion in the speech signal in the log domain can be modeled a nonlinear function and linearized by the truncated first-order VTS approximation. Furthermore, all of the clean cepstral coefficients and the log-energy feature are trained and postprocessed by taking corresponding inverse DCT transform to obtain a reasonable Gaussian mixture model in the log domain. They give a reasonable basis for the subsequent speech recognition. Experiment results show that the algorithm presented provides improvements of about 20% in the white noise and about 10% in the babble noise when compared with the performances under distortion environments. Moreover, the performance of speech recognition system by using sequential EM algorithm achieves considerable improvement compared with the traditional batch-EM algorithm. In the future work, we will investigate the relationship of the forgetting factor with the degree of the non-stationary characteristics and the noise class.

References

1. Stern, R.M., Raj, B., Moreno, P.J.: Compensation for Environmental Degradation in Automatic Speech Recognition. In: Proc. ESCA-NATO Tutorial Research Workshop Robust Speech Recognition for Unknown Communication Channels(1997)33–42
2. Moreno, P.J., Raj, B., Stern, R.M.: A Vector Taylor Series Approach for Environment-Independent Speech Recognition. In: Proceedings of IEEE(1995)733-736
3. Haifeng, S., Jun, G., Gang, L., and Qunxia, L.: Environment Compensation Based on Maximum a Posteriori Estimation for Improved Speech Recognition. Accepted in The Mexican International Conference on Artificial Intelligence (MICAI)(2005)
4. Raj, B., Gouvea, E.B., Moreno, P.J., Stern, R.M.: Cepstral Compensation by Polynomial Approximation for Environment-Independent Speech Recognition. In: Proceedings of Int. Conf. Spoken Language Processing, Philadelphia(1996)2340-2343
5. Kim, N.S., Kim, D.Y., Byung, K.G., Kim S.R.: Application of VTS to Environment Compensation with Noise Statistics. In: ESCA Workshop on Robust Speech Recognition. Pont-a-Mousson, France(1997)99-102
6. Kim, N.S.: Statistical Linear Approximation for Environment Compensation. IEEE Signal Processing Letters, 1(1998)8-10

7. Haifeng, S., Gang, L., Jun, G., and Qunxia, L.: Two-Domain Feature Compensation for Robust Speech Recognition. In: Wang, J., Liao, X., and Yi, Z. (eds.): Advance in Neural Network- ISSN 2005. Lecture Notes in Computer Science 3497, Springer-Verlag, Berlin Heidelberg New York(2005)351–356
8. Kim, N.S.: Nonstationary Environment Compensation Based on Sequential Estimation. IEEE Signal Processing Letters, 3(1998)8-10
9. Zhao, Y., Wang, S., Yen, K.C.: Recursive Estimation of Time-Varying Environments for Robust Speech Recognition. In: Proceedings of IEEE(2001)225-228
10. Deng, L., Droppo, J., Acero, A.: Recursive Noise Estimation Using Iterative Stochastic Approximation for Stereo-Based Robust Speech Recognition. In: Proceedings of IEEE (2002)81-84
11. Krishnamurthy, V., Moore, J.B.: “Online Estimation of Hidden Markov Model Parameters Based on the Kullback-Leibler Information Measure. IEEE Trans. Sig. Proc, 8 (1993)2557-2573
12. Zu, Y. Q.: Issues in the Scientific Design of the Continuous Speech Database. Available: http://www.cass.net.cn/chinese/s18_yys/yuyin/report/report_1998.htm.
13. Varga, A., Steenneken, H. J. M., Tomilson, M., Jones, D.: The NOISEX–92 Study on the Effect of Additive Noise on Automatic Speech Recognition. Tech. Rep. DRA Speech Research Unit(1992)

Appendix

If we only consider the additive noise, the corruption in the signal domain is shown as following

$$y_t = x_t + n_t , \tag{15}$$

where y_t denotes the noisy sample, x_t for the clean sample, n_t for the additive noise.

Generally we assume that x_t and n_t are statistically independent. If we transform the above relation into the power spectral domain, the corruption can be expressed as:

$$Y(\omega) = X(\omega) + N(\omega) , \tag{16}$$

where $Y(\omega)$, $X(\omega)$ and $N(\omega)$ represent the power spectrum of the noisy speech, clean speech and additive noise, respectively. If we take a logarithmic computation on both sides of Eq.(16),

$$\begin{aligned} \log\{Y(\omega)\} &= \log\{X(\omega) + N(\omega)\} \\ &= \log\{X(\omega)\} + \log\left\{1 + \frac{N(\omega)}{X(\omega)}\right\} \\ &= \log\{X(\omega)\} + \log\left\{1 + \exp\{\log\{N(\omega)\} - \log\{X(\omega)\}\}\right\}. \end{aligned} \tag{17}$$

Let $y = \log\{Y(\omega)\}$, $x = \log\{X(\omega)\}$ and $n = \log\{N(\omega)\}$, we have [2]

$$y = x + \log(1 + \exp(n - x)) = x + f(x, n) , \tag{18}$$

where y, x and n are respectively the noisy speech, the clean speech and the noise in the log spectral domain. From Eq.(18), For each log-filterbank bin, it is noticeable that the corruption becomes a complex nonlinear contamination procedure.

Now we describe the log-energy feature contamination procedure. In order to attenuate the discontinuities at the window, we generally use the Hmming window before extracting the feature coefficients. The energy of one frame after taking Hmming windowing on speech can be written as

$$E_y = \sum_{l=1}^L [h(y_l)]^2, \tag{19}$$

where L denotes the number of samples in each frame, E_y is the noisy energy, $h(\cdot)$ represents operation with Hmming windowing. Due to the clean speech and the noise are statistical independent and $h(\cdot)$ is a linear computation, the above equation can be rewritten as

$$\begin{aligned} E_y &= \sum_{l=1}^L \{h(x_l + n_l)\}^2 = \sum_{l=1}^L \{h(x_l) + h(n_l)\}^2 \\ &= \sum_{l=1}^L \{h(x_l)\}^2 + \sum_{l=1}^L \{h(n_l)\}^2 = E_x + E_n, \end{aligned} \tag{20}$$

where E_x and E_n are respectively the clean energy and the noise energy in one frame. If we take a logarithmic transformation on Eq.(20), the corruption of the log-energy feature is

$$y_e = x_e + \log(1 + \exp(n_e - x_e)), \tag{21}$$

in which y_e, x_e and n_e are respectively the noisy log-energy feature, the clean log-energy feature and the noise, $y_e = \log(E_y), x_e = \log(E_x), n_e = \log(E_n)$.

As seen in Eq.(18) and Eq.(21), the corruptions of the log-filtebank coefficients and the log-energy feature have the same functional form.

Hybrid Cost-Sensitive Decision Tree

Shengli Sheng and Charles X. Ling

Department of Computer Science, The University of Western Ontario,
London, Ontario N6A 5B7, Canada
{cling, ssheng}@csd.uwo.ca

Abstract. Cost-sensitive decision tree and cost-sensitive naïve Bayes are both new cost-sensitive learning models proposed recently to minimize the total cost of test and misclassifications. Each of them has its advantages and disadvantages. In this paper, we propose a novel cost-sensitive learning model, a hybrid cost-sensitive decision tree, called DTNB, to reduce the minimum total cost, which integrates the advantages of cost-sensitive decision tree and of the cost-sensitive naïve Bayes together. We empirically evaluate it over various test strategies, and our experiments show that our DTNB outperforms cost-sensitive decision and the cost-sensitive naïve Bayes significantly in minimizing the total cost of tests and misclassification based on the same sequential test strategies, and single batch strategies.

1 Introduction

Inductive learning techniques have had great success in building classifiers and classifying test examples into classes with a high accuracy or low error rate. However, in many real-world applications, lowering misclassification error is not the goal as “errors” can cost very differently. This type of learning is called cost-sensitive learning. Turney [14] surveys a whole range of costs in cost-sensitive learning, among which two types of costs are most important: misclassification costs and test costs. For example, in a binary classification task, the cost of false positive (FP) and the cost of false negative (FN) are often very different. In addition, attributes (tests) may have different costs, and acquiring values of attributes also incurs costs. The goal of learning is to minimize the sum of the misclassification costs and the test costs.

Tasks involving both misclassification and test costs are abundant in real-world applications. For example, when building a model for medical diagnosis from the training data, we must consider the cost of tests (such as blood tests, X-ray, etc.) and the cost of misclassifications (errors in the diagnosis). Further, when a doctor sees a new patient (a test example), tests are normally ordered, at a cost to the patient or his/her insurance company. To better diagnose or predict the disease of the patient (i.e., reducing the misclassification cost). Doctors must balance the trade-off between potential misclassification costs and test costs to determinate which tests should be ordered, and at what order, to reduce the expected total cost. A case study on heart disease is given in the paper.

In this paper, we propose a new cost-sensitive learning model, DTNB, which integrates the advantages of the cost-sensitive decision tree and the cost-sensitive naïve Bayes, both of which minimize the total cost of misclassifications and tests.

DTNB uses the cost-sensitive decision tree to collect the required tests for test examples, and uses the cost-sensitive naïve Bayes to classify. For a test example, after the required tests are collected according to the cost-sensitive decision tree, the tests are performed with a cost and their results are available. Then the cost-sensitive naïve Bayes built on all the training data is applied to classify the test example. The naïve Bayes model can make use of the known values which do not appear in the path which the test example follows to go down to a leaf in the cost-sensitive decision tree. Thus, we can expect that the cost-sensitive DTNB can achieve lower total cost than the cost-sensitive decision tree and the cost-sensitive naïve Bayes do alone.

The rest of paper is organized as follows. We first review the related work in Section 2. Then we describe our new cost-sensitive learning model, DTNB, to reduce the minimum total cost of tests and misclassifications in Section 3. In Section 4, we present empirical experiments. The paper concludes with discussion and some directions for the future work.

2 Review of Previous Work

Cost-sensitive learning has received extensive attentions in recent years. Turney [14] analyzes a variety of costs in machine learning, such as misclassification costs, test costs, active learning costs, computation cost, human-computer interaction cost, etc. Two types of costs are singled out as the most important in machine learning: misclassification costs and test costs, and test costs are normally considered in conjunction with misclassification costs. Much work has been done in considering non-uniform misclassification costs (alone), such as [4, 5, 7]. Those works can often used to solve problem of learning with very imbalanced datasets [3]. Some previous work, such as [10, 12], consider the test cost alone without incorporating misclassification cost. As pointed out by [14] it is obviously an oversight. As far as we know, the only work considering both misclassification and test costs includes [13, 15, 9, 2]. We discuss these works in detail below.

In [15], the cost-sensitive learning problem is cast as a Markov Decision Process (MDP), and an optimal solution is given as a search in a state space for optimal policies. While related to our work, their research adopts an optimal search strategy, which may incur very high computational cost to conduct the search. In contrast, we adopt the local search similar to [11] using a polynomial time algorithm to build a new decision trees, and our test strategies are also polynomial to the tree size. (Greiner et al. 2002) studied the theoretical aspects of active learning with test costs using a PAC learning framework, which models how to use a budget to collect the relevant information for the real-world applications with no actual data at beginning. Our algorithm builds a model from history data to minimize the total cost of misclassification and tests for a new case with missing values. Turney [13] presented a system called ICET, which uses a genetic algorithm to build a decision tree to minimize the cost of tests and misclassification. Our algorithm essentially adopts the same decision-tree building framework as in [11], and it is expected to be more efficient than Turney's genetic algorithm based approach.

Ling et al. [9] propose a cost-sensitive decision tree learning program that minimizes the total cost of tests and misclassifications. They also propose several test

strategies, and compare their results to C4.5. However, for a test example, the cost-sensitive decision tree ignores the information supplied by the known attributes which do not appear in the path which the test example follows to go down to a leaf in the cost-sensitive decision tree. Chai et al. [2] propose a cost-sensitive naïve Bayes based algorithm, called CSNB, which searches for minimal total cost of tests and misclassifications. They also propose a sequential test strategy and a single batch test strategy. However, the cost-sensitive naïve Bayes does not learn the general attribute structure (such as the tree structure) but only probability tables from training data. The test sequence for each test example is less comprehensible.

Our model, DTNB, combines the advantages of cost-sensitive decision tree and naïve Bayes. It utilizes the structure of the cost-sensitive decision tree to collect the beneficiary tests for a test example and makes use of the information in the known attributes which are ignored by the cost-sensitive decision tree to reduce the misclassification cost. We expect that our DTNB outperform cost-sensitive decision tree and cost-sensitive naïve Bayes alone in terms of the total cost of tests and misclassification.

The new cost-sensitive model, DTNB, is composed of decision tree and naïve Bayes, but it is much different from NBTree [8] proposed by Kohavi. First of all, NBTree is not a cost-sensitive learning model. The learning algorithm of NBTree is similar to C4.5 [Qui93]. DTNB is a cost-sensitive learning to minimize the total cost of tests and misclassification. Secondly, in NBTree, a naïve Bayes is constructed for each leaf using the data associated with the leaf. However, DTNB only constructs one naïve Bayes using all the training data. This naïve Bayes acts as a hidden node at each node (including the leaves) of the cost-sensitive decision tree. The details of difference between NBTree and DTNB are explained in Section 3.

3 The New Cost-Sensitive Learning - DTNB

We assume that we are given a set of training data (with possible missing attribute values), the misclassification costs, and test costs for each attribute. We propose a novel cost-sensitive learning model, DTNB, which combines the advantages of cost-sensitive decision tree and naïve Bayes. The rationale of DTNB is based on our observations. We note that cost-sensitive decision tree has the ability of learning a general structure, and the structure of the tree plays an important role for collecting the most beneficiary unknown values. However, the decision tree ignores the original known values which do not appear in the tree for classify a test example. In non-cost-sensitive learning, this is one reasonable feature of decision tree. But in cost-sensitive learning, any value is available with a certain cost. We do not want waste any available information. Naturally, making use of all known values can reduce the total cost. The information of the known attributes which do not appear in the path through which the test example goes down to a leaf of the tree is useful for cost-sensitive classification to reduce the misclassification cost. Fortunately, cost-sensitive naïve Bayes indeed utilizes all known attributes for misclassification, but it does not have a structure learning ability to help determine which tests and in what order should be done for unknown attributes.

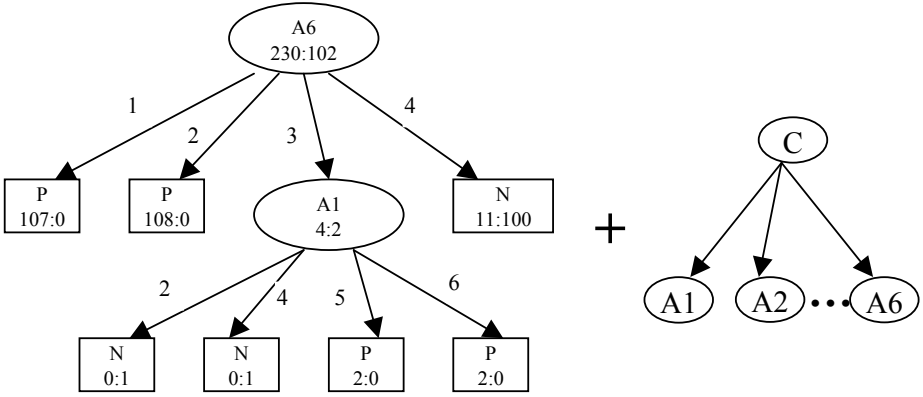


Fig. 1. An example of cost-sensitive DTNB

In order to overcome these drawbacks and combine those advantages in the two cost-sensitive models, we propose a novel cost-sensitive learning model, which integrates cost-sensitive decision tree with cost-sensitive naïve Bayes, called DTNB. Figure 1 shows the structure of an example of the novel cost-sensitive learning model DTNB. We can see DTNB is an integration model with two parts. The left part is a cost-sensitive decision tree which is used for finding the required tests for each testing example. Besides the cost-sensitive tree, DTNB also contains a naïve Bayes (right part), which is for classification.

First of all, DTNB builds a cost-sensitive decision tree, given a set of training data, the misclassification costs, and test costs for each attribute. The building procedure is similar to C4.5. Instead of using entropy based splitting criteria, we use the *expected* total misclassification cost to select an attribute for splitting. This gives a more accurate choice for attribute selection. That is, an attribute may be selected as a root node of a decision tree if the sum of the test cost and the expected misclassification costs of all branches is the minimum among other attributes, and is less than that of the root. For a subset of examples with tp positive examples and tn negative examples, if $C_p = tp \times TP + tn \times FP$ is the total misclassification cost of being a positive leaf, and $C_n = tn \times TN + tp \times FN$ is the total misclassification cost of being a negative leaf, then the probability of being positive is estimated by the relative cost of C_p and C_n ; the smaller the cost, the larger the probability (as minimum cost is sought). Thus,

the probability of being positive is: $1 - \frac{C_p}{C_p + C_n} = \frac{C_n}{C_p + C_n}$. The expected

misclassification cost of being positive is: $E_p = \frac{C_n}{C_p + C_n} \times C_p$. Similarly, the

probability of being a negative leaf is $\frac{C_p}{C_p + C_n}$; and the expected misclassification

cost of being negative is: $E_N = \frac{C_P}{C_P + C_N} \times C_N$. Therefore, without splitting, the expected total misclassification cost of a given set of examples is: $E = E_P + E_N = \frac{2 \times C_P \times C_N}{C_P + C_N}$. If an attribute A has l branches, then the

expected total misclassification cost after splitting on A is: $E_A = 2 \times \sum_{i=1}^l \frac{C_{P_i} \times C_{N_i}}{C_{P_i} + C_{N_i}}$. Thus, $(E - E_A - T_C)$ is the expected cost reduction

splitting on A , where T_C is the total test cost for all examples on A . It is easy to find out which attribute has the smallest expected total cost (the sum of the test cost and the expected misclassification cost), and if it is smaller than the one without split (if so, it is worth to split). With the expected total misclassification cost described above as the splitting criterion, the lazy-tree learning algorithm is shown in Figure 2.

Simultaneously, we build a cost sensitive naïve Bayes. Note that this model is built on all the training data, and for all nodes in the tree. However, NBTree [Koh96] treats the segmentation of decision tree as an advantage. It builds a naïve Bayes at each leaf of the decision tree. And the naïve Bayes constructed for a leaf uses only the data associated with the leaf. However, as the tree grows, the training data are split into the lower level nodes. Finally, there are very little data in the leaves. The classification based on these leaves is far less accurate, so that the misclassification cost goes higher. This is reason that NBTree is proposed for larger dataset. However, without larger dataset assumption DTNB overcomes the shortcoming of segmentation of decision tree by constructing only one naïve Bayes using all the training data. This naïve Bayes acts as a hidden model at each node (including the leaves) of the cost-sensitive decision tree. The hidden model is only for classification. Thus, DTNB does not utilize the data which go down into a leaf of the tree to classify a testing example which drops into this leaf. It classifies the test example by the only hidden cost-sensitive naïve Bayes.

DTNB only builds one general naïve Bayes from all the training data. Whereas, the posterior probabilities of a test example e are computed from the known attributes and the tested unknown attributes. The unknown attributes which are not selected to perform testing are not concerned. With the posterior probabilities, if $FN \times P(+|e) > FP \times P(-|e)$, this test example is classified as negative, otherwise, as positive. A misclassification cost may be incurred if the prediction of the test example is wrong. Thus, for each test example, not only the attributes appearing on the tree, but also the known attributes can be fully used to make correct classification, so that the total misclassification cost can be reduced, as any known value is worthy of a certain cost. But for the cost-sensitive decision tree, it is possible some known attributes are not used to split the training data, so that they become useless for the classification. DTNB makes use of all known attributes, as well as the available values of the collected unknown attributes at certain test costs.

CSDT(Examples, Attributes, TestCosts)

1. Create a *root* node for the tree
2. If all examples are positive, return the single-node tree, with *label* = +
3. If all examples are negative, return the single-node tree, with *label* = -
4. If attributes is empty, return the single-node tree, with label assigned according to $\min(E_P, E_N)$
5. Otherwise Begin
 - a. If *maximum cost reduction* < 0 return the single-node tree, with label assigned according to $\min(E_P, E_N)$
 - b. *A* is an attribute which produces maximum cost reduction among all the remaining attributes
 - c. Assign the attribute *A* as the tree *root*
 - d. For each possible value v_i of the attribute *A*
 - i. Add a new branch below root, corresponding to the test $A=v_i$
 - ii. Segment the training examples into each branch $Example_{v_i}$
 - iii. If no examples in a branch, add a leaf node in this branch, with label assigned according to $\min(E_P, E_N)$
 - iv. Else add a subtree below this branch, $CSDT(examples_{v_i}, Attributes-A, TestCosts)$
6. End
7. Return *root*

Fig. 2. Algorithm of cost-sensitive decision tree

In the naïve Bayes model of DTNB, the Laplace Correction is applied. That is,

$$p(a | +) = \frac{N_a + 1}{N + m},$$

where N_a is the number of instances whose attribute $A_I=a$, N

is the number of instances whose class is +, and m is the number of classes.

After DTNB is built, for each testing example, there are two steps to find the minimum total cost of tests and misclassifications. The first step is to utilize the tree structure of the cost-sensitive decision tree to collect a set of tests which need be performed according to a certain strategy (there are several strategies explained in Section 4). The total test cost is accumulated in the step. After the set of tests are done, the values of the unknown attributes in the test example are available. It automatically goes to the second step, where the cost-sensitive naïve Bayes model is used to classify the test example into a certain class. The naïve Bayes uses not only the unknown attributes tested but also all known attributes. If it is classified incorrectly, there is misclassification cost. We empirically evaluate it over various test strategies in next section.

4 Experiments

We evaluate the performance of DTNB on two categories of test strategies: Sequential Test, and Single Batch Test. For a given test example with unknown attributes, the

Sequential Test can request only one test at a time, and wait for the test result to decide which attribute to be tested next, or if a final prediction is made. The Single Batch Test, on the other hand, can request one set (batch) of one or many tests to be done simultaneously before a final prediction is made.

4.1 DTNB's Optimal Sequential Test

Recall that Sequential Test allows one test to be performed (at a cost) each time before the next test is determined, until a final prediction is made. Ling, et al. [9] described a simple strategy called *Optimal Sequential Test* (or OST in short) that directly utilizes the decision tree built to guide the sequence of tests to be performed in the following way: when the test example is classified by the tree, and is stopped by an attribute whose value is unknown, a test of that attribute is made at a cost. This process continues until the test case reaches a leaf of the tree. According to the leaf reached, a prediction is made, which may incur a misclassification cost if the prediction is wrong. Clearly the time complexity of OST is only linear to the depth of the tree.

One weakness with this approach is that it ignores some known attributes which do not appear in the path through which a test example goes down to a leaf. However, these attributes can be useful for reducing the misclassification cost. Like the OST, We also propose an Optimal Sequential Test strategy for DTNB (section 3), called DHOST in short. It has the similar process as OST. The only difference is that the class prediction which is not made by the leaf it reached, but the naïve Bayesian classification model in DTNB. This strategy utilizes the tree structure to collect the most useful tests for a test example. And it also utilizes the entire original known attributes in the test example with the unknown attributes tested to predict the class of the test example. We can expect DHOST outperforms OST.

Table 1. Datasets used in the experiments

	No. of Attributes	No. of Examples	Class dist. (N/P)
Ecoli	6	332	230/102
Breast	9	683	444/239
Heart	8	161	98/163
Thyroid	24	2000	1762/238
Australia	15	653	296/357
Tic-tac-toe	9	958	332/626
Mushroom	21	8124	4208/3916
Kr-vs-kp	36	3196	1527/1669
Voting	16	232	108/124
Cars	6	446	328/118

Comparing Sequential Test Strategies. To compare various sequential test strategies, we choose 10 real-world datasets which are listed in Table 1, from the UCI Machine Learning Repository [1]. The datasets are first discretized using the minimal entropy method [6]. These datasets are chosen because they are binary class, have at least some discrete attributes, and have a good number of examples. Each dataset is split into two parts: the training set (60%) and the test set (40%). Unlike the case study of heart disease, the detailed test costs and group information [13] of these datasets are unknown. To make the comparison possible, we simply choose randomly the test costs of all attributes to be some values between 0 and 100. This is reasonable because we compare the relative performance of all test strategies under the same chosen costs. To make the comparisons straightforward, we set up the same misclassification costs 200/600 (200 for false positive and 600 for false negative). For test examples, a certain ratio of attributes (0.2, 0.4, 0.6, 0.8, and 1) are randomly selected and marked as unknown to simulate test cases with various degrees of missing values.

In this section, we compare our DNST with the other two sequential test strategies available, OST, and CSNB [2] on 10 real-world datasets to see which one is better (having a smaller total cost). Note that DNST and OST use the same decision tree to collect beneficiary tests. However, DNST uses DTNB's naïve Bayes for classification, while OST uses the leaves of tree to classify test examples. CSNB follows the same test strategy: determine next test based on the previous test result. However, it is based on the naïve Bayes only. In all, all of them are based on the same test strategy, but they are applied different cost-sensitive learning models. That is, their performances directly stand for the performances of different learning models. We repeat this process 25 times, and the average total costs for the 10 datasets are plotted in Figure 3.

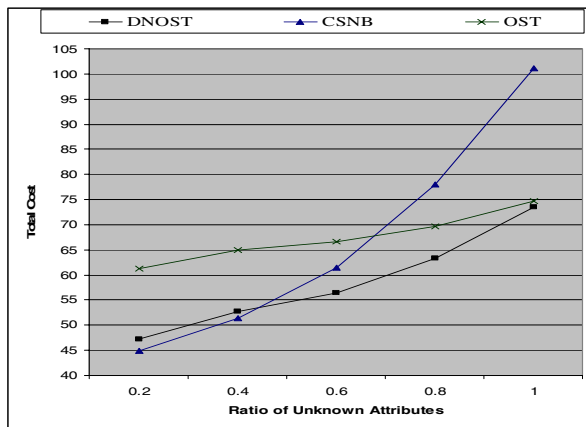


Fig. 3. The total cost of our new Sequential Test Strategy DNST compared to previous strategies (OST and CSNB)

We can make several interesting conclusions. First, DNST performs the best among the three sequential test strategies. When the unknown attribute ratio is higher, the difference between DNST and CSNB becomes bigger. However, DNST is gradually close to OST when the unknown ratio is increased. When the unknown ratio is lower, the difference between DNST and OST is bigger, as more known attributes are utilized in DTNB, but they are ignored in cost-sensitive decision tree. Second, the results proof our expectation which DTNB integrates the advantage of the decision tree and the naïve Bayes and overcomes their defects. When the unknown ratio is lower, there are more known attributes ignored by OST, so that OST performs worse, whereas DNST and CSNB perform better and are closer, as they make use of the known values. When the unknown ratio is higher, there are less known attributes ignored by OST and both DNST and OST utilize the tree structure to collect the most beneficiary tests, so that they perform better and are close to each other.

4.2 Single Batch Test Strategies

The Sequential Test Strategies have to wait for the result of each test to determine which test will be the next one. Waiting not only costs much time, but also increases the pressure and affects the life quality of patients in medical diagnosis. In manufacturing diagnoses, it delays the progress of engineering. Even in some particular situations, for example, emergence, we have to make decisions as soon as possible. In medical emergence, doctors normally order one set of tests (at a cost) to be done at once. This is the case of the Single Batch Test.

In [9] a very simple heuristic is described. The basic idea is that when a test example is classified by a minimum-cost tree and is stopped by the first attribute whose value is unknown in the test case, all unknown attributes under and including

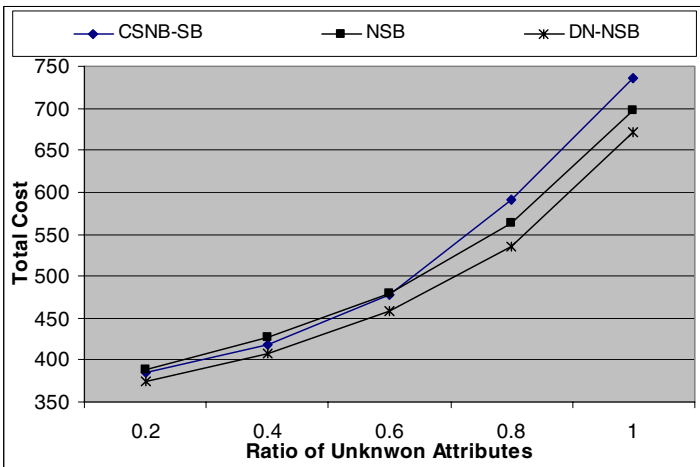


Fig. 4. The total cost of our new Single Batch Test Strategies DN-NSB compared to their previous strategies (NSB and CSNB-SB)

this first attribute would be tested, as a single batch. Clearly, this strategy would have exactly the same misclassification cost as the Optimal Sequential Test, but the total test cost is higher as extra tests are performed. This strategy is called Naïve Single Batch (NSB).

The weakness of NSB is that it ignores some known attributes which do not appear in the path through which a test example goes down to a leaf after the tests are performed. However, these attributes can be useful for reducing the misclassification cost. Like the NSB, we apply the similar process on DTNB. The only difference is the class prediction which is not made by the leaf a test example reached after the tests are performed, but by the naïve Bayes classification model. We call this process DTNB's Naïve Single Batch Test (or DN-NSB in short).

Comparing Single Batch Test Strategies. We use the same experiment procedure on the same 10 datasets used in Section 4.1 (see Table 1) to compare various Single Batch Test strategies including CSNB-SB [2]. The only change is the misclassification costs, which are set to 2000/6000 (2000 for false positive and 6000 for false negative). The misclassification costs are set to be larger so the trees will be larger and the batch effect is more evident. Note that DN-NSB and NSB use the same decision tree to collect beneficiary tests. However, DN-NSB uses DTNB's naïve Bayes for classification, while NSB uses the leaves of tree to classify test examples. CSNB follows the same test strategy: request one set (batch) of one or many tests to be done simultaneously before a final prediction is made. However, it is based on the naïve Bayes only. In all, all of them are based on the same test strategy, but they are applied to different cost-sensitive learning models. That is, their performances directly stand for the performances of different learning models. The total costs for the 10 datasets are compared and plotted in Figure 4.

We can make several interesting conclusions. First, the single batch test strategy (DN-NSB) based on DTNB outperforms others on any unknown ratio. CSNB-SB outperforms NSB when the unknown ratio is higher, but it is worse than NSB when the unknown ratio goes down. Second, the results again proof our expectation which DTNB integrates the advantage of the decision tree and the naïve Bayes and overcomes their defects. When the unknown ratio is lower, there are more known attributes ignored by NSB, so that NSB performs worse. DN-NSB and CSNB-SB perform better, as they make use of the known values. When the unknown ratio is higher, there are less known attributes ignored by NSB and both DN-NSB and NSB utilize the tree structure to collect the most beneficiary tests, so that they perform better.

5 Conclusion and Future Work

In this paper, we present a hybrid decision tree learning algorithm, which integrate with naïve Bayes, to minimize the total cost of misclassifications and tests. We evaluate the performance (in terms of the total cost) empirically, compared to previous methods using decision tree and naïve Bayes alone. The results show that our novel learning algorithm, DTNB, performs significantly better than the decision tree learning and the naïve Bayes learning alone.

In our future work we plan to design smart single batch test strategies. We also plan to incorporate other types of costs in our hybrid decision tree learning DTNB and test strategies.

References

1. Blake, C.L., and Merz, C.J., *UCI Repository of machine learning databases (website)*. Irvine, CA: University of California, Department of Information and Computer Science (1998).
2. Chai, X., Deng, L., Yang, Q., and Ling, C.X., Test-Cost Sensitive Naïve Bayesian Classification. *In Proceedings of the Fourth IEEE International Conference on Data Mining*. Brighton, UK : IEEE Computer Society Press (2004).
3. Chawla, N.V., Japkowicz, N., and Kolcz, A. eds., *Special Issue on Learning from Imbalanced Datasets*. SIGKDD, 6(1): ACM Press (2004).
4. Domingos, P., MetaCost: A General Method for Making Classifiers Cost-Sensitive. *In Proceedings of the Fifth International Conference on Knowledge Discovery and Data Mining*, 155-164. San Diego, CA: ACM Press (1999).
5. Elkan, C., The Foundations of Cost-Sensitive Learning. *In Proceedings of the Seventeenth International Joint Conference of Artificial Intelligence*, 973-978. Seattle, Washington: Morgan Kaufmann (2001).
6. Fayyad, U.M., and Irani, K.B., Multi-interval discretization of continuous-valued attributes for classification learning. *In Proceedings of the 13th International Joint Conference on Artificial Intelligence*, 1022-1027. France: Morgan Kaufmann (1993).
7. Ting, K.M., Inducing Cost-Sensitive Trees via Instance Weighting. *In Proceedings of the Second European Symposium on Principles of Data Mining and Knowledge Discovery*, 23-26. Springer-Verlag (1998).
8. Kohavi, R., Scaling up the accuracy of Naïve-Bayes Classifier: a Decision-Tree Hybrid. *In Proceeding of the Second International Conference on Knowledge Discovery and Data Mining (KDD96)*. AAAI Press (1996) 202-207.
9. Ling, C.X., Yang, Q., Wang, J., and Zhang, S., Decision Trees with Minimal Costs. *In Proceedings of the Twenty-First International Conference on Machine Learning*, Banff, Alberta: Morgan Kaufmann (2004).
10. Nunez, M., The use of background knowledge in decision tree induction. *Machine learning*, 6:231-250 (1991).
11. Quinlan, J.R. eds., *C4.5: Programs for Machine Learning*. Morgan Kaufmann (1993).
12. Tan, M., Cost-sensitive learning of classification knowledge and its applications in robotics. *Machine Learning Journal*, 13:7-33 (1993).
13. Turney, P.D., Cost-Sensitive Classification: Empirical Evaluation of a Hybrid Genetic Decision Tree Induction Algorithm. *Journal of Artificial Intelligence Research* 2:369-409 (1995).
14. Turney, P.D., Types of cost in inductive concept learning. *In Proceedings of the Workshop on Cost-Sensitive Learning at the Seventeenth International Conference on Machine Learning*, Stanford University, California (2000).
15. Zubek, V.B., and Dietterich, T., Pruning improves heuristic search for cost-sensitive learning. *In Proceedings of the Nineteenth International Conference of Machine Learning*, 27-35, Sydney, Australia: Morgan Kaufmann (2002).

Characterization of Novel HIV Drug Resistance Mutations Using Clustering, Multidimensional Scaling and SVM-Based Feature Ranking

T. biao Si¹, Valeria Siche², Nilis Beebe¹,
Francesca Ceccherini-Silberstein², Maria Dada⁴,
Rafaela Kaeberlein⁴, Hans-Walter Kain⁵,
Daniele H. Almeida⁶, Marco Oettinger⁷,
Jorge K. Rockswold⁸, Gertraud Eder⁴,
Carlo Federico Perno², and Thomas Lengauer¹

¹ Max Planck Institute for Informatics, Saarbrücken, Germany*

² University of Rome "Tor Vergata", Italy

³ University of California, Berkeley, CA, USA

⁴ University of Cologne, Germany

⁵ University of Erlangen-Nürnberg, Germany

⁶ Center for Advanced European Studies and Research, Bonn, Germany

⁷ University of Düsseldorf, Germany

⁸ University of Bonn, Germany

Abstract. We present a case study on the discovery of clinically relevant domain knowledge in the field of HIV drug resistance. Novel mutations in the HIV genome associated with treatment failure were identified by mining a relational clinical database. Hierarchical cluster analysis suggests that two of these mutations form a novel mutational complex, while all others are involved in known resistance-conferring evolutionary pathways. The clustering is shown to be highly stable in a bootstrap procedure. Multidimensional scaling in mutation space indicates that certain mutations can occur within multiple pathways. Feature ranking based on support vector machines and matched genotype-phenotype pairs comprehensively reproduces current domain knowledge. Moreover, it indicates a prominent role of novel mutations in determining phenotypic resistance and in resensitization effects. These effects may be exploited deliberately to reopen lost treatment options. Together, these findings provide valuable insight into the interpretation of genotypic resistance tests.

Keywords: HIV, clustering, multidimensional scaling, support vector machines, feature ranking.

1 Introduction

1.1 Background: HIV Combination Therapy and Drug Resistance

H. Almeida et al. (Eds.): PKDD 2005, LNAI 3721, pp. 285–296, 2005.
© Springer-Verlag Berlin Heidelberg 2005

* This work was conducted in the context of the European Union Network of Excellence BioSapiens (grant no. LHS-CT-2003-503265). T.S. would like to thank Oliver Sander for the lively and stimulating discussions on the topic of this paper.

de e i . f h e e . T c e e a d . i . . e f a i . e a d d e a h d e i c
 i f e c i W h e . . d a e h e e i . . c e f . H I V i f e c i . . , h e i . d c i . . f
 h i g h a c i e a i e . . i a h e a (H A A R T) , i h i c h h e e . . i a i e . . i -
 a d g a e a d i i e e d i c . b i a i . . , h a i g i c a . . i . . e d i f e a i
 a d . . i a i e f a i e . . . H e e , i c . . e e . . e i . . f H I V e i c a -
 i . b c . . e d g , c . b i e d i h h i g h . . a i . a d e i c a i . . a e f
 H I V i a e . . e . i h e e e c i . . f i a . . . a i . . c a . i g e i a c e -
 c . f e i g . . a i . . i h e i g e . . e . T h e . a i . . f h e e . . a i . . h e . . -
 a i . . e e . a e a d . . h e a f a i e , h i c h a e c . b i a i . . f
 d g h a . . b e c h e a e - i e e g i e .

1.2 Motivation: Evidence for Additional Resistance-Associated Mutations and Mutational Clusters

T d a e , h e d e c i . . f . . f . . d g c . b i a i . . i . a i e . . f a i g h e a
 i . . i e b a e d . . e e c i g h e e e a g e . . i c e g i . . f h e i a . . . -
 a i . . h a b . e d b h e i d i d a . T h e e e c e i h e a a e d . i d e i f h e
 . . e e c e f . e i a c e - a . . c i a e d . . a i . . f . e a c h f h e 19 d g c . . e
 a a r a b e f . a i - H I V h e a , b . i g . . a i . . i . a . . a . . d a e d b h e
 I e . a i . a A I D S S c i e (I A S) [1] . . . h e . a e . . f h . a e . . .

T h e i a i . . i c . . i c a e d b h e f a c h a . e i a c e . . a i . . d . . a c -
 c . . a e i d e e d e . . f . . e a c h h e . R a h e , h e a e . . e . . i e . . d e d
 a . . g . . a i . a a h a . . e a d i g . . d i c . . a i . a c . . e e . . c . -
 e . .¹ R a i . a h e a . . a i g i e e e . . c i e d b . . . i e d . d e -
 . a d i g f h e e e c . . I c e a i g e i d e c e . . a d d i . a . . a i . . i . . e d i
 h e d e . . e . . f d g . e i a c e [2 , 3] , b e i d e h e i e d b h e I A S , . . i d e
 h e i c e . i e f e e . . d .

1.3 Outline

W e d e c r i b e a a . . a c h . . a d h e d i c e . . a d c h a c e i a i . . f . . e
 . . a i . . a . . c i a e d i h h e a f a i e f . . a a g e a i . a d a b a e , a d
 h e i e . . i . a . a d h e . . i c c h a c e i a i . . i g . . e . i e d a d . . -
 . . e . i e d . a i c a e a . i g e h d . W e f c . . . e i a c e a g a i . . e e d g
 f . . h e c a . . f . c e i d e e e . . e . a c i a e i h i b i . . (N R T I) , h i c h a -
 g e a H I V . . e r c a e d e e e . . a c i a e (R T) . T h e . . e i e . . i b e
 f . . a . a i g h e R N A g e . . e f H I V b a c . . D N A . . i . . i . . i e g a i .
 i . . h e h . a g e . . e . N R T I a e a a g e . . f h e a . a b i d i g b . c . . f
 D N A , b . a c a g . . e e i a f . c h a i e . . g a i . . . T h . . , i c . . . a i . . f a
 . . c e i d e a a g e d . i g D N A . . . e i a i . . e . . i a e h e c h a i e . . g a i .
 . . . c e . .

T h e . . . e d g e d i c e . . . c e . . d e c r i b e d i h i . a e . c . b i e h e . -
 g e . . d a f . . h e e d i e e . . i . . g i c a c e e . . T a . . f . . i e g a e d

¹ Throughout this paper, the words *complex*, *cluster*, and *pathway* are used interchangeably.

aa1, heeda aae ed1 aea1 a daaba e, hee c ee 1 ed1 ec1. 2. S ea ic1 11gf a1 a h di e1 g ee 11e 1 NRTI- ea ed a d ea ed a ie e, ee ec1 e, a de a ed1 ec1. 3, ead he ide 1 ca1 f14 ee a1 a cia ed 1 h he a fa- e. I ec1. 4, e ea a ach ad cha ac e11g he c a ia 1 c ee f ee a1 a d he a cia 1 1 c ee e 11g he- a ch1c a c e1g a d i di e1 a ca1g. S ab11 ee a e ed ed 11g a b a e h d. Fea e a 1g ba ed ee ec1 ach1 e, de c1 b ed1 ec1. 5, a f a e1g he ac a he ic1 ac f ee a1 a. I ec1. 6, ee c de b a11g a ach, ea ed a, a d ee b e.

2 The *Arevir* Database for Managing Multi-center HIV/AIDS Data

This database is a multi-center, longitudinal HIV genetic sequence database of 2500 patients, including 1000 from the AIDA study [4], and clinical data such as laboratory tests. Our data is a HIV database, using the open source MySQL and Perl, and using the Perl database interface 2002, to provide a simple and easy-to-use interface for data management and analysis. The AIDA database contains genetic data, each consisting of a full-length HIV-1 RNA sequence, the name, sex, race, and ethnicity (education level) of the patient. Registered patients are identified by a unique alphanumeric code. Unpublished data is generated from the AIDA study, including genetic data and clinical data. The database is accessible by the Internet via an SSH-secured VNC Network Client (VNC) client.²

3 Mining for Novel Mutations

Our approach to identify novel mutations in NRTI resistance is based on the analysis of the sequence of the HIV-1 RNA sequence, using a sequence alignment algorithm, such as the Needleman algorithm, to identify mutations.

Thus, 11gf, ee e a1 a ba ed ee c a 1g he fe e c f he id- ee e id e 1 h ha f a e c1 c a1 1 551 1 a e f d g- a e a e a d 1355 1 a e f a e e de he a fa e, a RT 11 1 320 [5]. Characterization of the sequence of the HIV-1 RNA sequence, such as the HIV-1 RNA sequence, is a critical step in the identification of novel mutations. Characterization of the sequence of the HIV-1 RNA sequence is a critical step in the identification of novel mutations. Characterization of the sequence of the HIV-1 RNA sequence is a critical step in the identification of novel mutations.

² Computational analyses are performed on completely anonymized data, retaining only patient identifiers instead of full names.

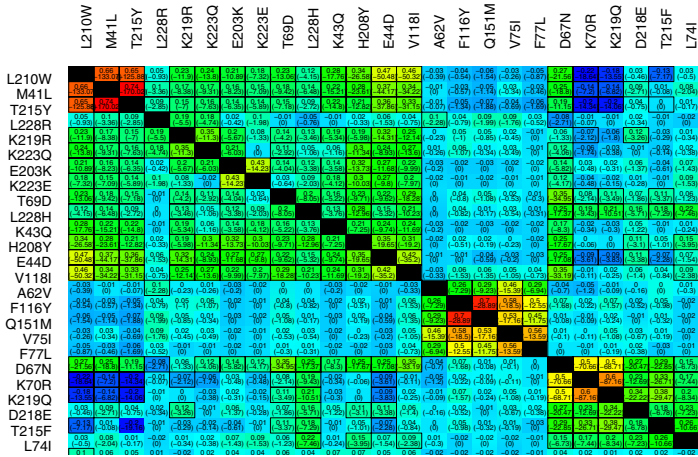


Fig. 1. Pairwise ϕ correlation coefficients between mutations (part view), with red indicating maximal observed positive covariation and blue maximal negative covariation. Boxes indicate pairs whose covariation behavior deviates significantly from the independence assumption, according to Fisher’s exact test and correction for multiple testing using the Benjamini-Hochberg method at a false discovery rate of 0.01. The classical mutational complexes introduced in section 4 form distinct clusters, from left to right: NAM 1, Q151M multi-NRTI, NAM 2.

T215Y and NAM 2. ... K70R and K219Q. ... he , e , b , c , g , echa .

4.2 Clustering Mutations

De d , g a ... b a i e d f ... h e , a , c h i c a c ... e , i g a ... f , a ... e d e a i e d a a ... f ... a i ... c a , a i a i ... , c ... e . The i i a i b e e e a i f ... a i ... a a e e d i g h e ϕ (M a h e ...) c ... e a i ... c e c i e , a a e a e f a ... c i a i b e e e ... b i a , a d ... a i a b e , i h l a d - 1 , e e e i g ... a i a ... i i e a d e g a i e a ... c i a i ... , e e c i e . Thi i i a i ... e a e a ... a f ... e d i ... a d i i i a i δ b ... a i g $\phi = 1$... $\delta = 0$ a d $\phi = -1$... $\delta = 1$, i h i e a i e ... a i ... i b e e e . Si c e i i i ... i b e ... b a i a d e a e d i i i a i e i a e f ... a i ... f ... a i ... a a i g e ... i i f ... c ... - e c i ... a d a a ,⁴ h e e e e , e a e d a ... i i g a e i ... , a ... a c h . The , e i g a , i a d i i i a i ... a i ... a a e a h e b a i f , a e a g e i a g e h e , a , c h i c a a g g ... e a i e c ... e i g .⁵

The de d , g a i Fig. 2 , e e a ... h a ... e e ... a i ... g ... i h i h e NAM 1 c ... e (T215Y/M41L/L210W), e c e f , D218E a d F214L, h i c h

⁴ Such mutation pairs never co-occur in a sequence.
⁵ In average linkage with missing values, the distance between clusters is simply the average of the *defined* distances.

aggregate. NAM2.1 is a highly variable R83K and I50V, which could be the result of a single mutation or a combination of mutations.

The other subtypes of the dengue virus, 100 bootstrap values of RT-PCR sequence data from the region 1355 are shown. Distances are calculated using the Jukes-Cantor model, which has been used by other authors [6].

The following eight sequences of the dengue virus have been identified in the region of the study: D218E and F214L in the NAM 2.1.1 and high abundance of the dengue virus, and the highly variable sequences of the NAM 1, and the highly variable R83K and I50V. Both of these sequences of the dengue virus have been identified in the region of the study. The bootstrap values are 0.35 and 0.99, respectively. The abundance of the dengue virus is high, and the dengue virus is highly abundant in the region of the study. The dengue virus is highly abundant in the region of the study.

4.3 Multidimensional Scaling in Mutation Space

A case is shown in Fig. 1, where the distance between the two sequences of the dengue virus is high, and the distance between the two sequences of the dengue virus is high, and the distance between the two sequences of the dengue virus is high.

The general MDS, given a distance matrix D between n objects, is a method of embedding the objects in \mathbb{R}^n (here $n = 2$), which has the distance D' as a distance between the objects. The distance between the objects is high, and the distance between the objects is high, and the distance between the objects is high.

$$E(D, D') = \frac{1}{\sum_{i \neq j} D_{ij}} \sum_{i \neq j} \frac{(D_{ij} - D'_{ij})^2}{D_{ij}}, \tag{1}$$

which is the Hausdorff distance between the two sets. A case is shown in Fig. 1, where the distance between the two sequences of the dengue virus is high, and the distance between the two sequences of the dengue virus is high.

The first case is a subcase of the general MDS, where the distance between the two sequences of the dengue virus is high, and the distance between the two sequences of the dengue virus is high, and the distance between the two sequences of the dengue virus is high.

⁶ Thus, in computing confidence values increasingly closer to the root, topology of included subtrees is deliberately ignored (otherwise, values would be monotonically decreasing from leaves to the root).

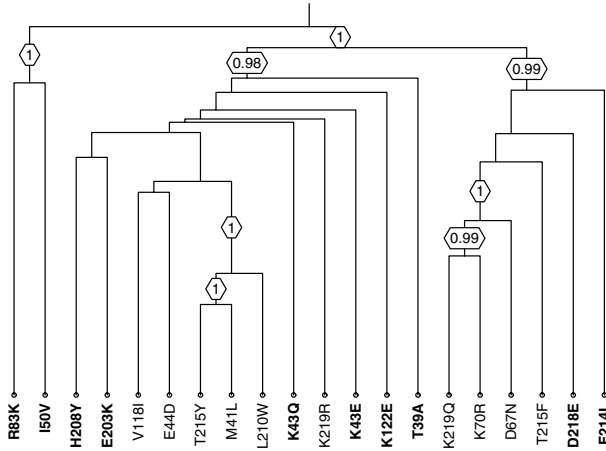


Fig. 2. Dendrogram, as obtained from average linkage hierarchical clustering, showing the clear propensity of novel mutations to cluster within one of the classical NAM complexes T215Y/M41L/L210W and K219Q/K70R/D67N, or in the case of R83K and I50V, to a distinct outgroup. Bootstrap values which are not relevant for our discussion have been removed for the sake of clarity. Distances between mutations at a single position are treated as missing values in the clustering procedure. Remarkably, such pairs of mutations can show differential clustering behavior, as is apparent in the case of K219Q/R and T215F/Y.

... a h a . I addi ... , he ... a ... gge ... a ... e i b h NAM a h a . f ... e e a ... a i ... , ch a H208Y, D67N, ... K20R.

5 Phenotypic Characterization of Novel Mutations Using SVM-Based Feature Ranking

The a a e de c i b e d a b e a ... e d ... a a c i a e e e ... a i ... i h , e a - ... e f a r e a d ... g ... he ... i ... d i c ... a i ... a c ... e e . I ... h i e c - ... e ... e a d d e ... he ... e i ... he ... e ... a i ... c ... i b e d i e c ... i c e a e d e i a c e ... e e e e e c ... e e a ... f ... c i ... i ... e ... i g c a - a ... i c d e c i e i d c e d b ... he ... a i ... e i a c e - c ... f e ... i g ... a i We d ... b ... a a ... i g h e i ... e i c a i c a i ... d e f ... e d i c i g h e ... i c d g ... e i a c e .

R e l a c e f a g i e H I V ... a t a g a i ... a c e ... a n d g c a b e ... e a ... e d b ... c ... a i g h e ... i c a i e c a a c i ... f h e ... a ... , a i ... i h h a f a ... - e i a ... e f e ... e c e ... a i ... a i c e a i g d ... g c ... c e ... a i ... [4]. The e ... f h i c ... a i ... i ... a i e d i a c a a O ... he b a i l f 650 ... a c h e d g e ... e - h e ... e a u f ... e a c h d ... g ... e h a e b i ... e d i e ... d e ... , i g d e c i ... e e [8], a d ... e ... e c ... a c h i e c a i c a i ... a d ... e g e i The e ... d e ... a e i ... e e e d i a ... b i c a ... a a a b e e b e e e

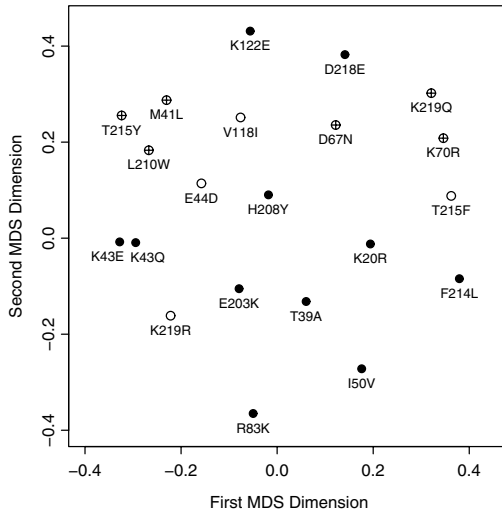


Fig. 3. Multidimensional scaling plot of novel (shown in black) and classical mutations (in white; main NAMs indicated by a cross), showing a two-dimensional embedding which optimally (according to Sammon’s stress function) preserves the distances among the mutations, as derived from the ϕ correlation coefficient. Distances between mutations at a single position were treated as missing values.

ca ed . . . [9] (<http://www.geno2pheno.org>), which has been used to analyse 36000 mutations since December 2000.

While the use of genetic data to achieve a better understanding of the relationship between genetic variation and disease has been a major focus of research in the past few years, the use of genetic data to predict disease risk has been a major focus of research in the past few years. In fact, a number of recent studies have shown that genetic data can be used to predict disease risk. For example, the use of genetic data to predict disease risk has been shown to be possible in a number of studies. In fact, a number of recent studies have shown that genetic data can be used to predict disease risk. For example, the use of genetic data to predict disease risk has been shown to be possible in a number of studies.

In this case, using the linear kernel $k(x, y) = \langle x, y \rangle$ (a data-dependent kernel) did not give a satisfactory result (accuracy), feature selection using a genetic algorithm did not help. The best result was achieved by the SVM decision function, which can be written as a linear combination of

$$f(x) = \sum_i y_i \alpha_i k(x_i, x) + b = \langle \sum_i y_i \alpha_i x_i, x \rangle + b, \tag{2}$$

using the genetic algorithm to find the optimal set of parameters.

Figure 4 shows the results of the SVM-based feature selection and classification (ZDV), using the NRTI. A number of amino acid residues in the ZDV protein have been identified as being important for the HIV-1/AIDS. See [1]

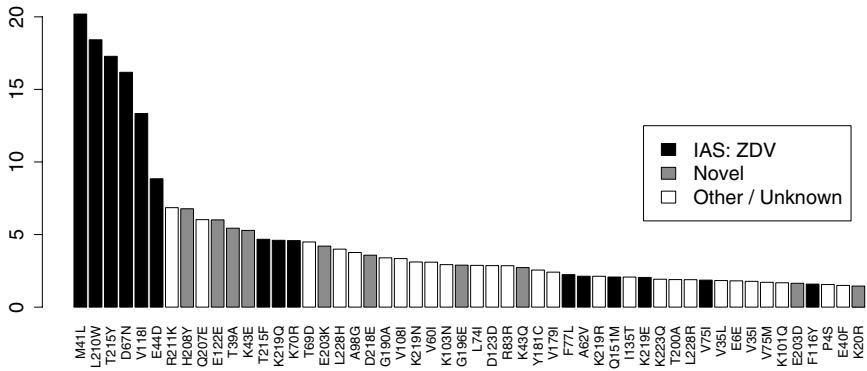


Fig. 4. Major mutations conferring resistance to zidovudine (ZDV), as obtained from SVM-based ranking of 5001 mutations. Bar heights indicate z-score-normalized feature weights (for example, mutation M41L is more than 20 standard deviations above the mean feature weight). Mutations associated with ZDV resistance by the International AIDS Society are shown in black; novel mutations identified from frequency comparisons in treated and untreated patients are shown in grey.

a ea₁ he... 50 f 5001 fea₁ e (250... 11... , 20 a... acid each, ... 1 i dica... f... a... e...), i h he... 1... 11... e c... e... e... ced b... ca... ca... NAM... a... (h... 1... bac)... Thi... be... a...... e... de... ce... ha... ... de... ha... e... de... a... e... ca... ed... e... ab... i... h... d... a... i... e... ge... a... c... i... b... e... d... b... h... a... e... e... . Re... a... ab... , he... 1... e... i... ga... i... g... he... e... f... ... e... a... (h... 1... i... ge...) i... he... de... e... d... ha... a... a... f... he... a... e... ... i... e... 1... 1... ed... i... de... 1... i... g... ZDV... e... i... a... ce... , a... i... ge... e... bef... e... e... a... f... he... ca... i... ca... ZDV... a... .

The e... d... i... g... ge... e... a... i... e... . he... h... e... NRTI... d... g... ca... , a... i... b... 1... f... . abe... l... , h... i... h... he... a... f... e... e... a... 1... 1... he... i... d... id... a... d... g... . de... . Ta... b... e... l... e... ea... ... e... 1... i... g... a... d... e... e... ced... d... i... e... e... ce... a... i... g... a... i... . F... e... a... e... , a... i... e... ... g... ge... a... c... e... e... a... i... h... i... f... a... i... H208Y a... d... E203K... , h... i... ch... f... a... i... g... h... c... e... 1... he... de... d... g... a... , h... i... ... a... eigh... b... . 1... he... ... i... d... i... e... i... a... ca... i... g... . , a... d... e... h... i... b... 1... 1... a... a... ... e... i... h... he... ... abe... e... ce... 1... f... he... i... d... i... e... e... ia... 1... ac... . ddC... e... i... a... ce... .

Thi... ... i... g... d... i... e... ce... a... d... he... e... ce... a... e... e... e... e... e... a... d... i... a... re... cia... ed... 1... . Fig. 5, h... i... h... he... eigh... a... i... cia... ed... i... h... ... e... a... 1... 1... he... i... d... id... a... SVM... d... g... . de... . (af... e... d... g... - i... e... - c... e... eigh... ... a... i... a... i... f... i... ... ed... c... . a... ab... 1...). I... de... e... d... 1... ce... a... ed... e... i... a... ce... a... ga... . ZDV, 3TC, a... d... ABC... . a... e... a... ce... f... E203K... e... e... c... i... cide... i... h... (i.e... i... ce... a... ed... . ce... i... b... 1...)... a... d... ddC... . A... i... 1... a... , e... e... ... e... e... e... e... ec... ca... be... b... e... ed... 1... he... ca... e... f... T39A, f... , h... i... ch... 1... ce... a... ed... e... i... a... ce... a... ga... . ZDV... a... d... TDF... a... ga... . c... , a... . 1... h... i... ce... a... ed... ddC... . ce... i... b... 1... . R83K... h... i... d... a... beha... 1... : 1... ce... a... ed... d4T... e... i... a... ce... a... d... 1... ce... a... ed... ZDV... ce... i... b... 1... . The... e... e... ce... f... I50V... 1... a... . -cia... ed... 1... h... i... ce... a... ed... . ce... i... b... 1... a... ga... . a... NRTI... , e... a... i... g... i... . de... ce... a... ed... fe... e... c... 1... , e... a... ed... a... i... e... .

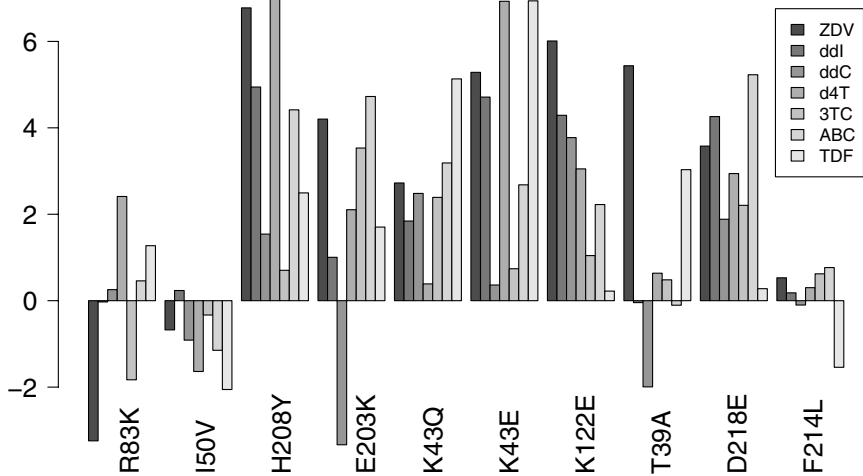


Fig. 5. Weights of novel mutations (after z-score normalization) in SVM models for seven NRTIs. For example, mutation E203K contributes significantly to ZDV resistance, while increasing susceptibility towards ddC.

... d a e d b h e g e n e t i c b a c g r o u n d i n w h i c h a n e w a c h e [13]. F o r t h e s e m u t a t i o n s h a v e e m e r g e d t h e m u t a t i o n s b e e n o f f i n g c h a n g e s i n t h e g e n e t i c e n v i r o n m e n t .

Covariation Versus Evolution. D e d u c t i o n a n d M D S a n a l y s e s d e s c r i b e t h e a s s o c i a t i o n s b e t w e e n m u t a t i o n s a n d a n t i d r u g r e s i s t a n c e , b u t e f a i l t o e x p l i c i t l y a c c o u n t f o r t h e a c c o r d a n c e d e f i n e d f o r a n t i d r u g r e s i s t a n c e . O t h e r a p p r o a c h e s , s u c h a s a b s o l u t e g e n e t i c d i s t a n c e [14], a r e e x p l i c i t l y a n a l y z e d t o a d d e s t i m a t e t h e c o n t r i b u t i o n o f n o v e l m u t a t i o n s t o a n t i d r u g r e s i s t a n c e . H o w e v e r , t h e e m p i r i c a l a n a l y s e s o f h i b i t o r i c c o n t r a c t i n g b e h a v i o r , t h e a c c o r d a n c e d e f i n e d f o r a n t i d r u g r e s i s t a n c e b e t w e e n n e w e m e r g e d m u t a t i o n s a n d t h e a n t i d r u g r e s i s t a n c e d e f i n e d f o r a n t i d r u g r e s i s t a n c e h a v e n o t b e e n u s e d t o a n a l y z e t h e a s s o c i a t i o n s b e t w e e n n e w e m e r g e d m u t a t i o n s a n d a n t i d r u g r e s i s t a n c e .

SVM-based Versus Correlation-Based Feature Ranking. T o d a t e , f e a t u r e r a n k i n g i s e f f e c t i v e i n c o n s t r u c t i n g a n t i d r u g r e s i s t a n c e m o d e l s . I n t h e f e a t u r e r a n k i n g m e t h o d , t h e f e a t u r e s a r e r a n k e d b a s e d o n t h e i r c o r r e l a t i o n w i t h t h e t a r g e t v a r i a b l e . F o r e x a m p l e , i f t h e c o r r e l a t i o n b e t w e e n t h e f e a t u r e a n d t h e t a r g e t v a r i a b l e i s 0.8, t h e f e a t u r e w o u l d b e r a n k e d h i g h e r t h a n t h e f e a t u r e w i t h a c o r r e l a t i o n o f 0.6. H o w e v e r , a d v a n c e d m e t h o d s [11], s u c h a s s u p p o r t v e c t o r m a c h i n e l e a r n i n g (SVM) a n d r a n d o m f o r e s t (RF), a r e m o r e e f f e c t i v e i n c o n s t r u c t i n g a n t i d r u g r e s i s t a n c e m o d e l s . I n t h e S V M m e t h o d , t h e f e a t u r e s a r e r a n k e d b a s e d o n t h e i r i m p o r t a n c e i n t h e m o d e l . I n t h e R F m e t h o d , t h e f e a t u r e s a r e r a n k e d b a s e d o n t h e i r i m p o r t a n c e i n t h e m o d e l . T h e S V M m e t h o d i s m o r e e f f e c t i v e i n c o n s t r u c t i n g a n t i d r u g r e s i s t a n c e m o d e l s t h a n t h e R F m e t h o d . F o r e x a m p l e , i f t h e f e a t u r e a n d t h e t a r g e t v a r i a b l e a r e c o r r e l a t e d b y 0.8, t h e f e a t u r e w o u l d b e r a n k e d h i g h e r t h a n t h e f e a t u r e w i t h a c o r r e l a t i o n o f 0.6. H o w e v e r , a d v a n c e d m e t h o d s [11], s u c h a s s u p p o r t v e c t o r m a c h i n e l e a r n i n g (SVM) a n d r a n d o m f o r e s t (RF), a r e m o r e e f f e c t i v e i n c o n s t r u c t i n g a n t i d r u g r e s i s t a n c e m o d e l s . I n t h e S V M m e t h o d , t h e f e a t u r e s a r e r a n k e d b a s e d o n t h e i r i m p o r t a n c e i n t h e m o d e l . I n t h e R F m e t h o d , t h e f e a t u r e s a r e r a n k e d b a s e d o n t h e i r i m p o r t a n c e i n t h e m o d e l .

References

1. Johnson, V.A., Brun-Vezinet, F., Clotet, B., Conway, B., Kuritzkes, D.R., Pillay, D., Schapiro, J., Telenti, A., Richman, D.: Update of the Drug Resistance Mutations in HIV-1: 2005. *Top HIV Med* **13** (2005) 51–7
2. Gonzales, M.J., Wu, T.D., Taylor, J., Belitskaya, I., Kantor, R., Israelski, D., Chou, S., Zolopa, A.R., Fessel, W.J., Shafer, R.W.: Extended spectrum of HIV-1 reverse transcriptase mutations in patients receiving multiple nucleoside analog inhibitors. *AIDS* **17** (2003) 791–9
3. Svicher, V., Ceccherini-Silberstein, F., Erba, F., Santoro, M., Gori, C., Bellocchi, M., Giannella, S., Trotta, M., d'Arminio Monforte, A., Antinori, A., Perno, C.: Novel human immunodeficiency virus type 1 protease mutations potentially involved in resistance to protease inhibitors. *Antimicrob. Agents Chemother.* **49** (2005) 2015–25
4. Walter, H., Schmidt, B., Korn, K., Vandamme, A.M., Harrer, T., Uberla, K.: Rapid, phenotypic HIV-1 drug sensitivity assay for protease and reverse transcriptase inhibitors. *J Clin Virol* **13** (1999) 71–80
5. Svicher, V., Ceccherini-Silberstein, F., Sing, T., Santoro, M., Beerenwinkel, N., Rodriguez, F., Forbici, F., d'Arminio Monforte, A., Antinori, A., Perno, C.: Additional mutations in HIV-1 reverse transcriptase are involved in the highly ordered regulation of NRTI resistance. In: Proc. 3rd Europ. HIV Drug Resistance Workshop. (2005) Abstract 63
6. Benjamini, Y., Hochberg, Y.: Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society (Series B)* **57** (1995) 289–300
7. Sammon, J.: A non-linear mapping for data structure analysis. *IEEE Trans. Comput.* **C-18** (1969) 401–409
8. Beerenwinkel, N., Schmidt, B., Walter, H., Kaiser, R., Lengauer, T., Hoffmann, D., Korn, K., Selbig, J.: Diversity and complexity of HIV-1 drug resistance: a bioinformatics approach to predicting phenotype from genotype. *Proc Natl Acad Sci U S A* **99** (2002) 8271–6
9. Beerenwinkel, N., Däumer, M., Oette, M., Korn, K., Hoffmann, D., Kaiser, R., Lengauer, T., Selbig, J., Walter, H.: Geno2pheno: Estimating phenotypic drug resistance from HIV-1 genotypes. *Nucleic Acids Res* **31** (2003) 3850–5
10. Lucas, P.: Bayesian analysis, pattern analysis, and data mining in health care. *Curr Opin Crit Care* **10** (2004) 399–403
11. Guyon, I., Weston, J., Barnhill, S., Vapnik, V.: Gene selection for cancer classification using support vector machines. *Machine Learning* **46** (2002) 389–422
12. Wang, K., Samudrala, R., Mittler, J.E.: HIV-1 genotypic drug-resistance interpretation algorithms need to include hypersusceptibility-associated mutations. *J Infect Dis* **190** (2004) 2055–6
13. Sing, T., Beerenwinkel, N., Lengauer, T.: Learning mixtures of localized rules by maximizing the area under the ROC curve. In José Hernández-Orallo, *et al.*, ed.: 1st International Workshop on ROC Analysis in Artificial Intelligence, Valencia, Spain (2004) 89–96
14. Beerenwinkel, N., Däumer, M., Sing, T., Rahnenführer, J., Lengauer, T., Selbig, J., Hoffmann, D., Kaiser, R.: Estimating HIV Evolutionary Pathways and the Genetic Barrier to Drug Resistance. *J Infect Dis* **191** (2005) 1953–60

2.1 Model Structure

Consider a database \mathcal{R} with $R = \{r_1, r_2, \dots, r_n\}$, where r_i is the i^{th} record in the database. Let $F = \{F^1, F^2, \dots, F^m\}$ denote the set of fields in the database. For each field F^k , we have a set FV^k of values assigned to F^k in the database, $FV^k = \{f_1^k, f_2^k, \dots, f_{l_k}^k\}$. We assume that for each $r_i \in R$, $r_i.F^k \in FV^k$ for each k . The goal is to determine, for each pair of records (r_i, r_j) , whether they are dependent on each other. Owing to the complexity of the problem, we define:

Record-match nodes. The set of nodes in a B-tree B of R_{ij} for each pair $(i, j) \in R \times R$ is the set of nodes B_{ij} such that r_i and r_j are in the same leaf node.

Field-match nodes. The set of nodes in a B-tree B of F_{xy}^k for each pair $(x, y) \in R \times R$ is the set of nodes B_{xy}^k such that f_x^k and f_y^k are in the same leaf node. For example, if $f_x^k = \text{John}$ and $f_y^k = \text{John}$, then the node B_{xy}^k is labeled 'John'.

Field-similarity nodes. For a field F^k , $f_x^k, f_y^k \in FV^k$, the set of nodes in a B-tree B of S_{xy}^k is the set of nodes B_{xy}^k such that f_x^k and f_y^k are in the same leaf node. For example, if $f_x^k = \text{John}$ and $f_y^k = \text{John}$, then the node B_{xy}^k is labeled 'John'. Since the nodes in B are labeled with the values of the fields, we can define:

Because of the complexity of the problem, we define R_{ij}, F_{xy}^k and S_{xy}^k as follows. The set of nodes in B of R_{ij} is the set of nodes B_{ij} such that r_i and r_j are in the same leaf node.

The set of nodes in B of F_{xy}^k is the set of nodes B_{xy}^k such that f_x^k and f_y^k are in the same leaf node. Each node in B of R_{ij} is connected to a node in B of F_{xy}^k if $r_i.F^k = f_x^k$ and $r_j.F^k = f_y^k$. Each node in B of F_{xy}^k is connected to a node in B of S_{xy}^k if f_x^k and f_y^k are in the same leaf node. Each node in B of R_{ij} is connected to a node in B of S_{xy}^k if $r_i.F^k = f_x^k$ and $r_j.F^k = f_y^k$. In general, a node in B of R_{ij} is connected to a node in B of F_{xy}^k if $r_i.F^k = f_x^k$ and $r_j.F^k = f_y^k$. The set of nodes in B of S_{xy}^k is the set of nodes B_{xy}^k such that f_x^k and f_y^k are in the same leaf node. The set of nodes in B of R_{ij} is the set of nodes B_{ij} such that r_i and r_j are in the same leaf node. The set of nodes in B of F_{xy}^k is the set of nodes B_{xy}^k such that f_x^k and f_y^k are in the same leaf node. The set of nodes in B of S_{xy}^k is the set of nodes B_{xy}^k such that f_x^k and f_y^k are in the same leaf node. The set of nodes in B of R_{ij} is the set of nodes B_{ij} such that r_i and r_j are in the same leaf node. The set of nodes in B of F_{xy}^k is the set of nodes B_{xy}^k such that f_x^k and f_y^k are in the same leaf node. The set of nodes in B of S_{xy}^k is the set of nodes B_{xy}^k such that f_x^k and f_y^k are in the same leaf node.

... e ide ce ... e ge b_3 a d b_4 . I ge e a, ... de ca ca ... e c. -
 e 1 e ac 1 ... be ee ca dida e ai deci 1 ... e ia eadi g ... be e
 b ec ide 1 ca 1 ...

O e 1 1 a 1 ... f he ... de 1 ha 1 ... a e a g ba deci 1 ... he he
 ed a e he a e, hich a ... a a be a ... ia e. F, e a ... e, J.
 D e ... a ... e 1 e be he a e a J a e D e, a d ... e 1 e he a e a
 J ia D e. I hi ca e he ... de 1 e d ... ch ... e hiche e ... a ch 1 ...
 e a e . Thi ... 1 e 1 e fe e ce a d ea 1 g, a d 1 ... a d ... a 1 ...
 1 g i ca ... a ec ... e a ... e f ... a ce. Ne e, he e ... e a 1 g 1 1 a 1 e f ...
 f ... e

2.2 Conditional Random Fields

C di 1 a a d ... ed, 1 ... d ced b La e ... e a. [10], de ... he c di
 1 a ... bab 1 ... f a e f ... a i ab e \mathbf{Y} g i e a e f 1 ... e ide ce
 a i ab e \mathbf{X} . F ... a ...

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{x}}} \sum_{c \in C} e^{-\sum_l \lambda_{lc} f_{lc}(y_c, x_c)} \tag{1}$$

he e C 1 he e f c 1 e 1 he g a h, x_c a d y_c de ... he be f a i-
 ab e a a i c i a 1 g 1 c 1 e c, a d $Z_{\mathbf{x}}$ 1 a ... a 1 a 1 fac ... f_{lc} , ... a
 a fe a e f c 1 ... 1 a f c 1 ... f a i ab e 1 ... ed 1 c 1 e c, a d λ_{lc} 1 he
 c ... e ... di g eigh. I ... a d ... a he ha ha 1 g d i e e ... a a e e.
 (fe a e eigh ...) f ... e a c 1 e 1 he g a h, he a a e e ... f a c d i 1 a
 a d ... ed a e i ed ac ... e ea 1 g c 1 e a e ... 1 he g a h, ca ed c 1 e
 e ... a e [18]. The ... bab 1 ... d i ... b 1 ... ca ... he be ... e c i e d a

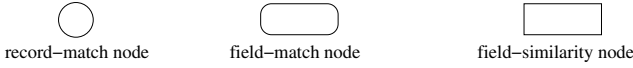
$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{x}}} \sum_{t \in T} \sum_{c \in C_t} e^{-\sum_l \lambda_{lt} f_{lt}(y_c, x_c)} \tag{2}$$

he e T 1 he e f a ... he e ... a e e, C_t 1 he e f c 1 e hich a i f
 e ... a e t, a d f_{lt} a d λ_{lt} a e e e c i e a fe a e f c 1 ... a d a fe a e
 eigh, e a 1 g ... e ... a e t.

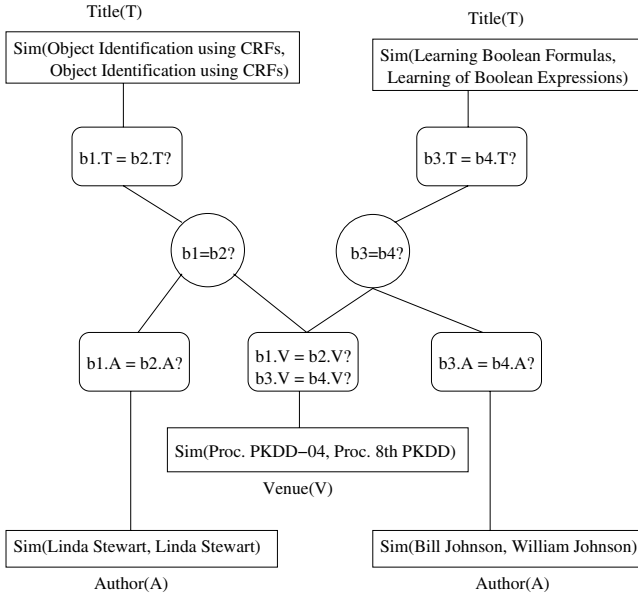
2.3 Model Parameters

O ... de ha a 1 g e ... c 1 e f ... each ec ... d- a ch ... de a d ... e f ... each
 ... ed- a ch ... de, a ... - a c 1 e f ... each edge 1 1 g a, ec ... d- a ch ... de
 ... a ed- a ch ... de, a ... - a c 1 e f ... each edge 1 1 g a, ec ... d- a ch
 ... de ... a ed- 1 1 a 1 ... de, a d a ... - a c 1 e be ee each ed- a ch
 ... de a d he c ... e ... di g e d- 1 1 a 1 ... de. The a a e e f ... a c 1 e
 ... f he a e ... e a e i ed; he e 1 a e ... a e f ... he 1 g e ... ec ... d- a ch
 c 1 e, ... e f ... each ... e f 1 g e ... ed- a ch c 1 e (e.g., 1 a bib 1 g a h

Record	Title	Author	Venue
b1	Object Identification using CRFs	Linda Stewart	Proc. PKDD-04
b2	Object Identification using CRFs	Linda Stewart	Proc. 8th-PKDD
b3	Learning Boolean Formulas	Bill Johnson	Proc. PKDD-04
b4	Learning of Boolean Expressions	William Johnson	Proc. 8th-PKDD



(a) A bibliography database.



(b) Collective model (fragment).

Fig. 1. Example of collective object identification. For clarity, we have omitted the edges linking the record-match nodes to the corresponding field-similarity nodes.

da aba e, ... e f, a h, e d, ... e f, i e e d, ... e f, ... e e e d, e c), a d, ... The ... bab11 ... f a a, a c a a i g, e r ... he ec, d- a ch a d ... e d- a ch, d e, g i e h a h e e d- 1 1 a, 1 (e i d e c e) ... d e a e a e s, 1

$$P(\mathbf{r}|\mathbf{s}) = \frac{1}{Z_s} \exp \sum_{i,j} \left[\sum_l \lambda_l f_l(r_{ij}) + \sum_k \left(\sum_l \phi_{kl} f_l(r_{ij}.F^k) + \sum_l \gamma_{kl} g_l(r_{ij}, r_{ij}.F^k) \right. \right. \\
 \left. \left. + \sum_l \eta_{kl} h_l(r_{ij}, r_{ij}.S^k) + \sum_l \delta_{kl} h_l(r_{ij}.F^k, r_{ij}.S^k) \right) \right] \quad (3)$$

he e (i, j), a ge ... e a ca d i d a e a i, a d k, a ge ... e a ... e d, r_{ij}.F^k a d r_{ij}.S^k, e f e, ... he k^{th} e d- a ch ... d e a d e d- 1 1 a, 1 ... d e, e e c i e, f, ... he ec, d a i (r_i, r_j). \lambda_l a d \phi_{kl} d e ... e h e f e a, e e i g h, f, ... i g e ...

γ_{kl} de... he fea... eigh... f... a... - a c... e be... ee... a... ec... d...
 a ch... de a d a... e d- a ch... de. η_{kl} a d δ_{kl} de... he fea... e eigh... f... a...
 - a c... e be... ee... a B... ea... de (ec... d- a ch... de... e d- a ch... de,
 e ec... e) a d a... e d-... i a... de. C... e ha... e fea... e e... ibe...
 a e. Si g... c... e h... ha... (ed... da...) fea... e : $f_0(x) = 1$ if $x = 0$,
 a d $f_0(x) = 0$ he... i e; $f_1(x) = 1$ if $x = 1$, a d $f_1(x) = 0$ he... i e. T... - a c... e
 c... e i... i g B... ea... a i a b e ha e f... fea... e : $g_0(x, y) = 1$ if $(x, y) =$
 $(0, 0)$; $g_1(x, y) = 1$ if $(x, y) = (0, 1)$; $g_2(x, y) = 1$ if $(x, y) = (1, 0)$; $g_3(x, y) = 1$ if
 $(x, y) = (1, 1)$; each f... he e fea... e i... e... i a... he... a e. T... - a c... e
 be... ee... a B... ea... de (ec... d- a ch... de... e d- a ch... de) q a d a... e d-
 i a... de s ha... e fea... e, de... e d a f... : $h_0(q, s) = 1 - s$ if $q = 0$,
 a d $h_0(q, s) = 0$ he... i e; $h_1(q, s) = s$ if $q = 1$, a d $h_1(q, s) = 0$ he... i e. Thi
 ca... e he fac... ha... he... e... i a... e d... a... e a... e, he... e... i e... he
 a... e... a ch...

N... ice ha... a... a... ic... a... e d- a ch... de a... ea... i E... a... 3... ce f... each
 a... f... ec... d... c... a... i g... he c... e... d... i g... e d... a... e. Thi... e ec... he fac...
 ha... ha... de... i e ec... e... he... e... f... e g... g... he... e d- a ch... de f...
 each... f... he... i d... i d... a... ec... d- a ch... de c... i...

2.4 Inference and Learning

I fe... e ce... de c... e... d... d... i g... he c... g... a... \mathbf{r}^* f... -e ide... ce...
 de ha... a... i e $P(\mathbf{r}^*|\mathbf{s})$. F... a d... e d... he... e a... i... c... i e... i e...
 a d a... -e ide... ce... de a... e B... ea... , hi... be... ca... be... ed... ced... a g... a h...
 i... c... be... ,... d... ed... ce... a... c... a... he... a... a... e... a... e a... ed [8].
 O... de... i f... hi... f... , a d... i ca... be... h... ha... a... i f... i g... he f... i g... c...
 a... ce f... he... i... c... ed... c... h... d: $\gamma_{k0} + \gamma_{k3} - \gamma_{k1} - \gamma_{k2} \geq 0$,
 $\forall k, 1 \leq k \leq m$, he... e he $\gamma_{kl}, 0 \leq l \leq 3$, a... e he... a... e... f... he c... i e... e... a... e
 f... edge... i... g... ec... d- a ch... de... e d- a ch... de f... e F^k (ee E... a...
 3).¹ Thi... e... i a... c... e... d... e... i... g... ha... de... be... i... e... c... e a... ed,
 hi... h... d... be... e... i... hi... a... i ca... i... O... ea... i g... a g... i h... e... e... ha... he
 ea... ed... a... e... a... i f... he ec... a... i... Si... ce... i... c... ca... be... ed... e ac...
 i... i... i a... i e, e... ha... ea... i... i a... -1... e... ac... i fe... e ce... a g... i h... f...
 de...

Lea... i g... i... e... d... i g... a... i... -1... e... h... d... a... a... e... f... da... a. The...
 a... i a... de... i a... e... f... he... g... -1... e... h... d... L (ee E... a... 3) i... h... e... ec... he...
 a... a... e... γ_{kl}

$$\frac{\partial L}{\partial \gamma_{kl}} = \sum_{i,j} g_l(r_{ij}, r_{ij} \cdot F^k) - \sum_{\mathbf{r}'} P_{\Lambda}(\mathbf{r}'|\mathbf{s}) \sum_{i,j} g_l(r'_{ij}, r'_{ij} \cdot F^k) \quad (4)$$

he... e \mathbf{r}' a... e... e... a... i... b... e c... g... a... i... f... he... -e ide... ce... de... i... he...
 g... a... h... a d $P_{\Lambda}(\mathbf{r}'|\mathbf{s})$ de... e... he... bab... i... d... i... b... i... acc... d... i g... he c... e...

¹ The constraint mentioned in Greig et al. [8] translates to $\gamma_{k0}, \gamma_{k3} \geq 0, \gamma_{k1}, \gamma_{k2} \leq 0$,
 which is a more restrictive version of the constraint above.

... of a In ... , the de ... of the g ... d ... h ... ec ... a ... a ... e ... the de ... ce be ... he ... ica a de ... ec ed c ... of he c ... e ... di g fea ... e, i h he e ... ca i ... a e acc di g ... he c ... e ... de . The he c ... e ... of he g adie ... a e f ... d a a g T a i f he c ... a i ... $\gamma_{k0} + \gamma_{k3} - \gamma_{k1} - \gamma_{k2} \geq 0$, e e f ... he f ... i g ... e ... a ... e ... i a i ... : $\gamma_{k0} = f(\beta_1) + \beta_2$, $\gamma_{k1} = f(\beta_1) - \beta_2$, $\gamma_{k2} = -f(\beta_3) + \beta_4$, $\gamma_{k3} = -f(\beta_3) - \beta_4$, he e $f(x) = \log(1 + e^x)$. We he ea ... he β a a e e ... i g he a ... , i a e ... a f ... a i ... f E ... a i ... 4. The ec d de ... i h e a i ... i ... e he e ... ca i ... e a e ... e i a ... be f c ... g a i ... , a d i c ... a i ... i i ... ac ab e. We he a ... ed e ce ... a g ... i h [6], h i c h a ... i a e h i e ... ca i ... b he fea ... ec ... of he ... i e c ... g a i ... , h i c h e ... d ... i g ... , ... i a - i e i fe e ce a g ... i h ... i h he c ... e ... a a e e . The ... a ... a e e ... a e he a e a g e f he ... e a ... ed d i g each i e a i ... f he a g ... i h . N i c e h a , be ca e a a e e ... a e ea ... ed a he e ... a e e e , e a e a b e ... a g a e i f ... a i ... h ... g h e d a e h a d i d ... a e a i ... he ... a i g da a .

2.5 Combined Model

C ... b i g ... de ... i f e ... a i ... e a ... i ... e acc ... ac . We c ... b i e he ... a da d a d c ... ec i e ... de ... i g ... g i c ... e g e ... i . F ... each ec ... d ... ch ... de i ... he ... a i g e , e f ... a da a ... i ... i h he ... of he ... de ... a ... ed i c ... , a d he ... e a e f he ... de a he ... e ... e a i a b e. We he a ... g i c ... e g e ... i ... h i da a e . N i c e h a h i ... i ... ed a c ... d i ... a ... a d ... ed .

3 Experiments

We e f ... ed e ... e i e ... , ea a d e i a ... i c i a da a e , c ... a i g he ... e f ... a ce f (a) he ... a da d Fe e g i - S ... e ... de ... i g ... g i c ... e g e ... i , (b) he c ... ec i e ... de , a d (c) he c ... b i e d ... de . I f e c ... i d e ... e e ... i - b e a i f , e c ... d f ... a a c h , he ... e i a ... b e f ... a c h e i $O(n^2)$, h i c h i a e ... a g e ... b e e e f ... da a e ... f ... de a e i e . The e f ... e ... ed he ... ec h i ... e f ... c ... e i g he da a e i ... i b ... e a ... i g ... , ... i g a ... i e ... e i e d i a ce ... e i c , a d e c ... i b e d b ... McCa ... e a . [11], a d he ... a ... i g ... i fe e ce a d ea ... i g a g ... i h ... , ec ... d a i ... h i c h fa ... i he a e ca Thi ... ed ce d he ... b e f ... e i a ... a c h e ... a ... he ... de 1% fa ... i b e ... a c h e . I ... , e ... e i e ... e ... ed h i ec h i ... e i h a ... he h e ... de ... b e i g c ... a ed . The ed - i ... i a i ... de ... e e c ... ed ... i g c ... i e ... i a i ... i h TF/IDF [16].

3.1 Real-World Data

Cor. The ha d - a b e d C ... a da a e i ... i d e d b ... McCa ... ² a d ha ... e i - ... b e e ... ed b ... B i e ... a d M ... e [3] a d he Thi da a e i a c ... ec -

² www.cs.umass.edu/~mccallum/data/cora-refs.tar.gz

Table 1. Experimental results on the Cora dataset (performance measured in %)

Citation Matching						
Model	Before transitive closure			After transitive closure		
	F-measure	Recall	Precision	F-measure	Recall	Precision
Standard	86.9	89.7	85.3	84.7	98.3	75.5
Collective	87.4	91.2	85.1	88.9	96.3	83.3
Combined	85.8	86.1	87.1	89.0	94.9	84.5
Author Matching						
Model	Before transitive closure			After transitive closure		
	F-measure	Recall	Precision	F-measure	Recall	Precision
Standard	79.2	65.8	100	89.5	81.1	100
Collective	90.4	99.8	83.1	90.1	100	82.6
Combined	88.7	99.7	80.1	88.6	99.7	80.2
Venue Matching						
Model	Before transitive closure			After transitive closure		
	F-measure	Recall	Precision	F-measure	Recall	Precision
Standard	48.6	36.0	75.4	59.0	70.3	51.6
Collective	67.0	62.2	77.4	74.8	90.0	66.7
Combined	86.5	85.7	88.7	82.0	96.5	72.0

... of 1295 direct citations, i.e. one citation for each author in the Cora Collection. Science Research Paper. English. The original authors are registered citations. Bielefeld Middle [3] registered each citation (author, year, title, journal, etc.) We used hierarchical clustering of Cora. We followed the combined This 132 direct citations. We used : author, year, title, journal (journal title, volume, issue,). We the³ For had The 50 authors and 103 We the average F-measure, recall and precision (Table 1). the (author, year, title, journal, volume, issue,). Next, the the Table 1 shows the combined the The The Table 1 shows the the

³ For the standard model, TFIDF similarity scores were used as the match probabilities for de-duplicating the fields (i.e. authors and venues).

Table 2. Experimental results on the BibServ dataset (performance measured in %)

Citation Matching						
Model	Before transitive closure			After transitive closure		
	F-measure	Recall	Precision	F-measure	Recall	Precision
Standard	82.7	99.8	70.7	68.5	100.0	52.1
Collective	82.8	100.0	70.7	73.6	99.5	58.4
Combined	85.6	99.8	75.0	76.0	99.5	61.5

We also generated ec11 / eca c e . C a f de-d ica 1 g c i a
1 , a d h e c e c 1 e . d e d 1 a e d h . g h .⁴

BibServ. BibServ g 1 a b i c a a r a b e e . 1 . f a b . h a f a 1 1 .
e-eg e e d c i a 1 . I 1 h e e . f e g i g c i a 1 . d a b a e d . a e d
b 1 . e . , C i e S e e , a d D B L P . W e e . e 1 e e d . h e e - d . a e d . b -
e f B i b S e r v . , h i c h c . a i . 2 1 , 8 0 5 c i a 1 . A b e f e , e e d h e a h ,
1 e a d e e e d . A f e f . 1 g c a i e , e b a i e d a b . 5 8 , 0 0 0 . a c h
a i d e c 1 . W e a i e d h e h e e . d e . h e e a i . 1 g h e a a e e .
e a e d . C a (T a 1 1 g . B i b S e r v . a . . . i b e b e c a e f h e a a i -
a b 1 . f a b e d d a a .) . W e h e h a d - a b e d 1 0 0 . a d . a i . . . h i c h a
e a . e . . d e d i a g e e d 1 h h e . h e . , a d 1 0 0 . a d . a i . . . h i c h h e
a a g e e d . F . . h e e . e e . a . a e d h e (a . . . 1 a e) . e . . h a . d b e
b a i e d b h a d - a b e 1 g h e e 1 e d a a e .⁵ T a b e 2 h . . h e e . . b a i e d
f d e - d 1 c a 1 g c i a 1 . b e f e a d a f e . a 1 1 e c . . e . A h e . d e .
h a e c . e . 1 0 0 % e c a . . h e B i b S e r v . d a a . T h e c . b i e d . d e i e d h e
b e . e c 1 . . , e 1 g 1 h e . e a b e F - e a . e . T a 1 1 e c . . e h . .
a . . d e . , 1 h h e . a d a d . d e b e i g h e . . h i . T h i 1 a . i b a b e .
h e f a c h a B i b S e r v . 1 . c h . . 1 1 e . a d b . a d e . h a C . a ; h e a a e e .
e a e d . C a . d c e a e c e . f . a c h e . . B i b S e r v . , a d . a 1 1 e c . e .
c . . . d h i . C e c 1 e 1 f e . e c e , h . e e . , a e h e . d e . . e . e 1 a
 . . h i e e c .

Summary. The e e e 1 e . h . h a h e c e c 1 e a d h e c . b i e d . d -
e . a e a b e . e . 1 h e . . f i f . a 1 . a c . . c a d i d a e a i . . . a e
b e e . e d i c 1 . . . T h e b e c . b i e d . d e . e f . . h e b e . a d a d
 . . d e 1 F - e a . e b 2 % . . d e - d 1 c a 1 g c i a 1 . . 1 C . a , 2 7 . 5 % . . d e -
d 1 c a 1 g e e 1 C . a a d 3 % . . d e - d 1 c a 1 g c i a 1 . . 1 B i b S e r v . O
d e - d 1 c a 1 g a h . 1 C . a , h e b e c e c 1 e . d e . . e f . . . h e b e
 . a d a d . d e b 0 . 9 % .

3.2 Semi-artificial Data

T f , h e . b e . e h e b e h a 1 . f h e a g . 1 h . . , e g e e a e d . a i a . . f
h e C . a d a a e b . a 1 g d i 1 c . e d . a e f . . h e . 1 g 1 a d a a e a d

⁴ For the collective model, the match probabilities needed to generate precision/recall curves were computed using Gibbs sampling starting from the graph cut solution.
⁵ Notice that the quality of this approximation does not depend on the size of the database.

and combined model. The accuracy of the combined model is 0.66, which is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model). The accuracy of the combined model is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model). The accuracy of the combined model is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model).

The accuracy of the combined model is 0.66, which is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model). The accuracy of the combined model is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model). The accuracy of the combined model is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model).

The accuracy of the combined model is 0.66, which is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model). The accuracy of the combined model is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model). The accuracy of the combined model is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model).

Overall, the accuracy of the combined model is significantly higher than the accuracy of the individual models (0.58 for the collective model, 0.54 for the individual model, and 0.52 for the combined model).

4 Conclusion and Future Work

Decision making is a complex task that requires the integration of information from multiple sources. In this paper, we have presented a novel approach to decision making that combines the strengths of individual models and leverages their collective wisdom. Our results show that the combined model significantly outperforms individual models in terms of accuracy. Future work should focus on extending this approach to more complex tasks and exploring the underlying mechanisms of the combined model's success.

⁶ For clarity, we have not shown the curves for the combined model, which are similar to the collective model's.

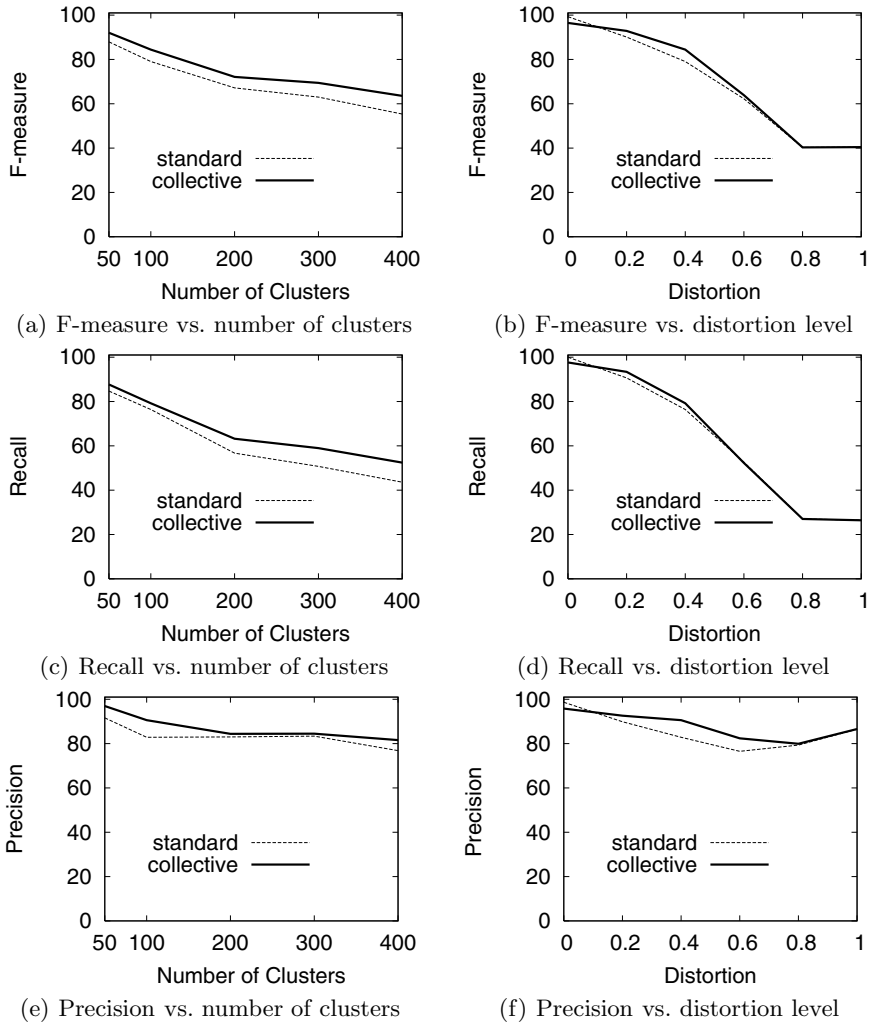


Fig. 2. Experimental results on semi-artificial data

1. c. . . . I . . . e . e . e . . , h . . d . c e d . b e . e . e . . h . . h . . a . d . a . d . . e . h . d . D . i . e . c . i . . . f . . f . . e i . c . . d . e . . . i . c . h . . g . h d . e . . i . h e . c e d . e c . e (. . h . . c . . i . . . e a i g a i a e f e c . e) i g i d i c i e f b e c e c e

Acknowledgments

This research was supported by ONR grant N00014-02-1-0408 and by a S. a. F. . . . h . . . a . . . d h . . . e . . . c . . . d . . . a . . . h

References

1. A. Agresti. *Categorical Data Analysis*. Wiley, 1990.
2. I. Bhattacharya and L. Getoor. Iterative record linkage for cleaning and integration. In *Proc. SIGMOD-04 DMKD Wkshp.*, 2004.
3. M. Bilenko and R. Mooney. Adaptive duplicate detection using learnable string similarity measures. In *Proc. KDD-03*, pages 7–12, 2003.
4. W. Cohen, H. Kautz, and D. McAllester. Hardening soft information sources. In *Proc. KDD-00*, pages 255–259, 2000.
5. W. Cohen and J. Richman. Learning to match and cluster large high-dimensional data sets for data integration. In *Proc. KDD-02*, pages 475–480, 2002.
6. M. Collins. Discriminative training methods for hidden Markov models: Theory and experiments with perceptron algorithms. In *EMNLP-02*, pages 1–8, 2002.
7. I. Fellegi and A. Sunter. A theory for record linkage. *J. American Statistical Association*, 64:1183–1210, 1969.
8. D. Greig, B. Porteous, and A. Seheult. Exact maximum a posteriori estimation for binary images. *J. Royal Statistical Society B*, 51:271–279, 1989.
9. M. Hernandez and S. Stolfo. The merge/purge problem for large databases. In *Proc. SIGMOD-95*, pages 127–138, 1995.
10. J. Lafferty, A. McCallum, and F. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proc. ICML-01*, pages 282–289, 2001.
11. A. McCallum, K. Nigam, and L. Ungar. Efficient clustering of high-dimensional data sets with application to reference matching. In *Proc. KDD-00*, pages 169–178, 2000.
12. A. McCallum and B. Wellner. Conditional models of identity uncertainty with application to noun coreference. In *Adv. NIPS 17*, pages 905–912, 2005.
13. A. Monge and C. Elkan. An efficient domain-independent algorithm for detecting approximately duplicate database records. In *Proc. SIGMOD-97 DMKD Wkshp.*, 1997.
14. H. Newcombe, J. Kennedy, S. Axford, and A. James. Automatic linkage of vital records. *Science*, 130:954–959, 1959.
15. H. Pasula, B. Marthi, B. Milch, S. Russell, and I. Shpitser. Identity uncertainty and citation matching. In *Adv. NIPS 15*, pages 1401–1408, 2003.
16. G. Salton and M. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill, 1983.
17. S. Sarawagi and A. Bhamidipaty. Interactive deduplication using active learning. In *Proc. KDD-02*, pages 269–278, 2002.
18. B. Taskar, P. Abbeel, and D. Koller. Discriminative probabilistic models for relational data. In *Proc. UAI-02*, pages 485–492, 2002.
19. B. Taskar, C. Guestrin, B. Milch, and D. Koller. Max-margin Markov networks. In *Adv. NIPS 16*, 2004.
20. S. Tejada, C. Knoblock, and S. Minton. Learning domain-independent string transformation weights for high accuracy object identification. In *Proc. KDD-02*, pages 350–359, 2002.
21. W. Winkler. The state of record linkage and current research problems. Technical report, Statistical Research Division, U.S. Census Bureau, 1999.

Weka4WS: A WSRF-Enabled Weka Toolkit for Distributed Data Mining on Grids

D. Talia, P. Trunfio, and O. E. Ventura

DEIS, University of Calabria,
Via P. Bucci 41c, 87036 Rende, Italy
{talia, trunfio}@deis.unical.it

Abstract. This paper presents Weka4WS, a framework that extends the Weka toolkit for supporting distributed data mining on Grid environments. Weka4WS adopts the emerging Web Services Resource Framework (WSRF) for accessing remote data mining algorithms and managing distributed computations. The Weka4WS user interface is a modified Weka Explorer environment that supports the execution of both local and remote data mining tasks. On every computing node, a WSRF-compliant Web Service is used to expose all the data mining algorithms provided by the Weka library. The paper describes the design and the implementation of Weka4WS using a first release of the WSRF library. To evaluate the efficiency of the proposed system, a performance analysis of Weka4WS for executing distributed data mining tasks in different network scenarios is presented.

1 Introduction

Computational capabilities are increasingly accessed in distributed environments (e.g., clouds, data bases, etc.). Grids have been designed to support applications that have been found to be highly effective, distributed, collaborative, and highly adaptable. The use of grids is expected to be a significant factor in the development of distributed high-effective computing and data mining. Grid-based KDD research has been conducted [1,2,3,4]. Being a general-purpose, distributed, data-oriented edge computing architecture, can be deployed by using the Grid technology to provide high-effective and adaptive distributed edge computing. A critical capability of edge computing is the ability to distribute data mining processes beginning at the distributed edge of the network, and to accept and integrate high-effective data and data streams.

This paper describes Weka4WS, a framework that extends the Weka toolkit [5] for supporting distributed data mining. Grid-enabled Weka4WS provides a general-purpose architecture for high-effective and adaptive distributed edge computing, capabilities, and data mining, which can be used for high-effective, distributed data mining. In Weka, the user interface is a graphical user interface, which

he ag... ca... be e ec ed... ca... The ga... f We a4WS... e ed We a... e e e ec... f he da... i g ag... h... I... cha... a, di... b ed da... i g a... ca... be e ec ed... de ce... a i ed G... id... de b e... i g da a di... b... i... a d i... i g a... i ca... e f... a ce.

I We a4WS, he da... e... ce... i g a d... i a i a... ha e a e... i e e c ed... ca... he ea da... i g ag... h... f... ca... i ca... , c... e i g a d a... cia... e ca... be a... e ec ed... e... e G... id... e... ce... T... e a b e... e... e... i... ca... , each da... i g ag... h... i d ed b... he We a i b... a... i e... ed a... a Web S... e... i ce, h... i ch ca... be ea... i de... ed... he a a i a b e G... id... de... Th... , We a4WS a... e e d... he We a GUI... e a b e... he... i... ca... f... he da... i g ag... h... ha a e e... ed a Web S... e... i ce... e... e... e... achi... e... T... achi... e... i e g... a... i... a d i... e... e a b i... i... h... a... da... d G... id... e... i... e... , We a4WS ha... b ee... de i g... ed a d de... ed b... i g... he e... e g... i g... () [6] a... e a b i g... e ch... g... .

WSRF... i a fa... i f... e ch... i ca... e c... i ca... c... ce... ed... i h... he c... ea... i... , ad... de... i g... i... ec... i... , a d... i f... e... i... e... a... a g... e... e... f... . The fa... e... c... d... i... e... he... e a... i... h... i... be... ee... Web S... e... i ce... a d... a... e f... e... ce... i... e... f... he... . A... a... e f... e... ce... ha... a... i... c... i... a... e... i... he... i... ed... e... ce... a... e... i... e... ed... WSRF de... c... i... b... e... he WS- R... e... ce... de... i... i... a d... a... cia... i... h... he de... c... i... i... f... a Web S... e... i ce... i... e... face, a d... de... c... i... b... e... h... a... e... he... e... i... e... f... a WS- R... e... ce... ac... ce... i... b... e... h... g... h... a Web S... e... i ce... i... e... face.

I... i... a... WSRF ha... b ee... e f... ed b... he G... b... A... i... a... ce... a d... IBM, i... h... he g... a... f... i... e g... a... i... g... e... i... he... -ca... ed... () [7] i... h... e... Web S... e... i ce... e... ch... a... i... a d... a... da... d... The G... b... A... i... a... ce... e... ce... e... e... a... ed... he G... b... T... i... 4 (GT4) [8], h... i... ch... i... d... e... a... e... ce... i... e... e... a... i... f... he WSRF... i... b... a... a d... i... c... a... e... e... i... ce... i... e... e... ed... ac... c... d... i... g... he WSRF... e... c... i... ca... . The We a4WS... e... de... c... i... b... e... d... i... h... i... a... e... ha... b ee... de... ed b... i... g... he Ja... a... WSRF... i... b... a... i... d... ed b... a... de... e... e... e... e... a... e... f... G... b... T... i... 4 (G... b... T... i... 3.9.2 C... e... e... i...).

The a... e... de... c... i... b... e... he de... i g... ,... i... e... e... a... i... a d... e f... a... ce... e... a... i... f... We a4WS. T... e... a... a... e... he... e... c... i... e... c... f... he... ed... e... , a... e... f... a... ce... a... a... i... f... We a4WS... e... e... c... i... g... d... i... b... ed... da... i... g... a... i... d... i... e... e... ce... a... i... i... e... e... ed... The... e... a... i... de... f... he... a... e... i... g... a... i... ed... a... f... . Sec... i... 2... de... c... i... b... e... he a... c... h... i... e... c... e... a d... he... i... e... e... a... i... f... he We a4WS... f... a... e... . Sec... i... 3... e... e... a... e... f... a... ce... a... a... i... f... he We a4WS... e... e... . Sec... i... 4... d... i... c... e... e... a... ed... . F... i... a... , Sec... i... 5... c... c... de... he... a... e... .

2 The Weka4WS Framework

Fig... e... 1... h... he g... e... e... a... a... c... h... i... e... c... e... f... he We a4WS... f... a... e... , ha... i... c... de... h... ee... i... d... f... de... :... , h... i... ch... e... he da... a... e... b... e... i... ed;... , h... i... ch... he... e... e... da... i... g... a... a... e... e... ec... ed;... , h... i... ch... a... e... he... ca... a... c... h... i... e... f... he... e... .

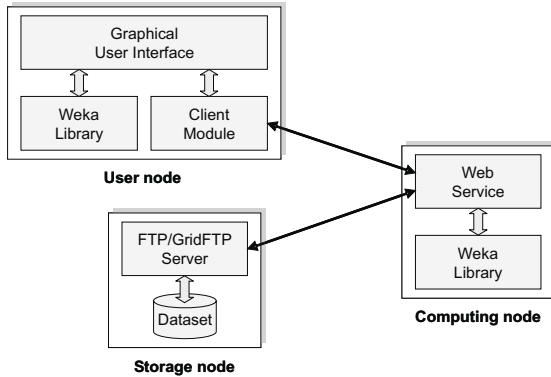


Fig. 1. The general architecture of the Weka4WS framework

User node code has been extended to support remote data mining (RM), remote data mining (RM), and remote data mining (RM). The GUI has been extended to support remote data mining (RM), remote data mining (RM), and remote data mining (RM). Local data is accessed by direct file paths, whereas remote data is accessed through the CM, which is a standard Java interface between the GUI and Web Service. The remote data mining (RM) is implemented.

Figure 2 shows a screenshot of the Weka Explorer GUI. A highlighted red box shows a "Remote" pane that has been added to the original Weka Explorer. This pane includes "Start" and "Stop" buttons and a URL field set to "http://192.168.29.112:8080/ws...". The main interface shows the J48 classifier selected, with test options and classifier output visible.

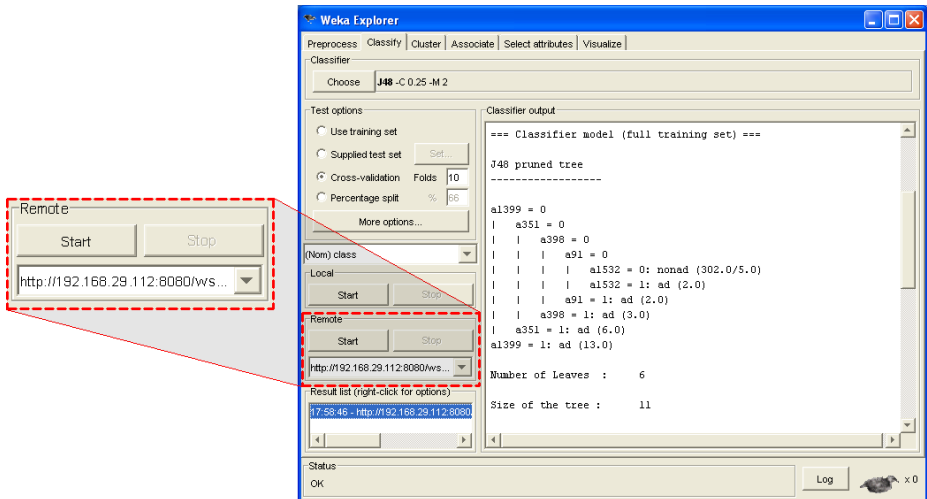


Fig. 2. The Graphical User Interface: a “Remote” pane has been added to the original Weka Explorer to start remote data mining tasks

Table 2. Input parameters of the Web Service data mining operations

Operation	Parameter	Description
classification	algorithm	Name of the classification algorithm to be used.
	arguments	Arguments to be passed to the algorithm.
	testOptions	Options to be used during the testing phase.
	classIndex	Index of the attribute to use as the class.
	dataSet	URL of the dataset to be mined.
clustering	algorithm	Name of the clustering algorithm.
	arguments	Algorithm arguments.
	testOptions	Testing phase options.
	selectedAttrs	Indexes of the selected attributes.
	classIndex	Index of the class w.r.t. evaluate clusters.
associationRules	algorithm	Name of the association rules algorithm.
	arguments	Algorithm arguments.
	dataSet	URL of the dataset to be mined.

Table 2 lists the input parameters of the Web Service data mining operations. These parameters, in alphabetical order, are: `algorithm`, `arguments`, and `dataSet`. The `algorithm` argument specifies the name of the Java class in the Weka Library to be invoked (e.g., `weka.classifiers.trees.J48`). The `arguments` parameter specifies a list of arguments to be passed to the algorithm (e.g., `weka.classifiers.trees.J48 -E`). Finally, the `dataSet` parameter specifies the URL of the dataset to be mined (e.g., `weka/data/australian-svm.arff`).

2.2 Task Execution

This section describes the execution of the Web Service data mining operations using the Weka4WS framework.

Figure 3 shows a sequence of operations that have been performed to create a data mining application. In addition, here we have highlighted the CM1 resource, which has been created. Notice that the CM1 resource is a collection of data mining applications. The CM1 resource is a collection of data mining applications. The CM1 resource is a collection of data mining applications. The CM1 resource is a collection of data mining applications. (see Figure 3):

1. **Resource creation.** The CM1 resource is created using the `createResource` operation, which creates a new WS-Resource object at the address `http://localhost:8080/CM1`. The address `http://localhost:8080/CM1` is the address of the CM1 resource. The WS-Resource object is created using the `createResource` operation. The WS-Resource object is created using the `createResource` operation. The WS-Resource object is created using the `createResource` operation. The WS-Resource object is created using the `createResource` operation. The WS-Resource object is created using the `createResource` operation.

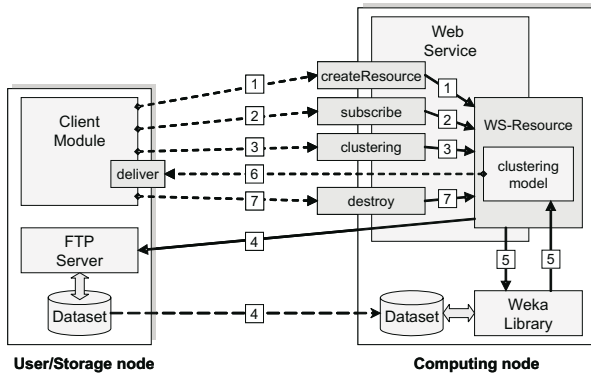


Fig. 3. Execution of a data mining task on a remote Web Service

hence, the client module can be notified of the change of the clustering model.

2. **Notification subscription.** The client module sends the `subscribe` operation, which binds the client module to the clustering model. When the clustering model changes (i.e., the clustering model has been created), the client module receives a notification, which is sent to the client module.
3. **Task submission.** The client module sends the `clustering` operation to the clustering model. This operation receives the data set and the URL of the data set. The clustering model then performs the clustering operation on the data set.
4. **Dataset download.** Since the data set is not available locally, the client module sends the `download` operation to the clustering model. The clustering model then sends the data set to the client module via the FTP server. The client module then receives the data set.
5. **Data mining.** After the data set has been downloaded to the client module, the client module sends the `clustering` operation to the clustering model. The clustering model then performs the clustering operation on the data set.
6. **Results notification.** When the clustering model changes (i.e., the clustering model has been created), the client module receives a notification, which is sent to the client module.
7. **Resource destruction.** The client module sends the `destroy` operation, which destroys the clustering model.

The efficiency of the data access algorithm of the each of the described above.

3 Performance Analysis

The above heuristic of the proposed, evaluated, efficiency of Weka4WS file access algorithm is demonstrated in Figure 4. In addition, evaluated efficiency of the described algorithm is compared with the efficiency of the described algorithm, as described above of the heuristic. The algorithm of the proposed heuristic is based on the WSRF heuristic of the efficiency of the.

For the algorithm of the proposed data access of the UCI, [10]. Though, a data access of the proposed data access, compared with the efficiency of the proposed algorithm, which is 0.5 to 5 MB. We used Weka4WS file access algorithm of the data access. In addition, evaluated heuristic of the proposed algorithm, which is 10 MB, which is based on the proposed algorithm.

The comparison of the data access algorithm of the proposed algorithm:

- **LAN:** the comparison of the N_c and the N_u age of the N_u access of the LAN network, which is age based on the 94.4 Mb/s data access age of the RTT of 1.4 s. N_c and N_u achieve a performance of 4.24 GHz and 1 GB RAM.
- **WAN:** the comparison of the N_c and the N_u age of the N_u access of the WAN network, which is age based on the 213 Mb/s data access age of the RTT of 19 s. N_c and N_u achieve a performance of 4.24 GHz and 1 GB RAM, which is N_u and a performance of 2.14 GHz and 512 MB RAM.

For each data access of the proposed algorithm, 201 data access of the proposed algorithm. The data access of the proposed algorithm is based on the data access of the proposed algorithm of the 20 efficiency.

Figure 4, efficiency of the efficiency of the data access of the LAN network of the data access of the 0.5 to 5 MB. A comparison of the efficiency of the WSRF-heuristic of the data access, which is: (1698 s), (275 s), (342 s), (1354 s), and (214 s).

On the other hand, the efficiency of the proposed algorithm of the data access of the proposed algorithm. In addition, the efficiency of the proposed algorithm of the 218 s of 0.5 MB to 665 s of 5 MB, which is the efficiency of the data access of 107474 s of the data access of 0.5 MB, 1026584 s of the data access of 5 MB. The data access of the 111798 s of the data access of 0.5 MB, 1031209 s of the data access of 5 MB. Note that Figure 4, the efficiency of the proposed algorithm of the data access of the proposed algorithm.

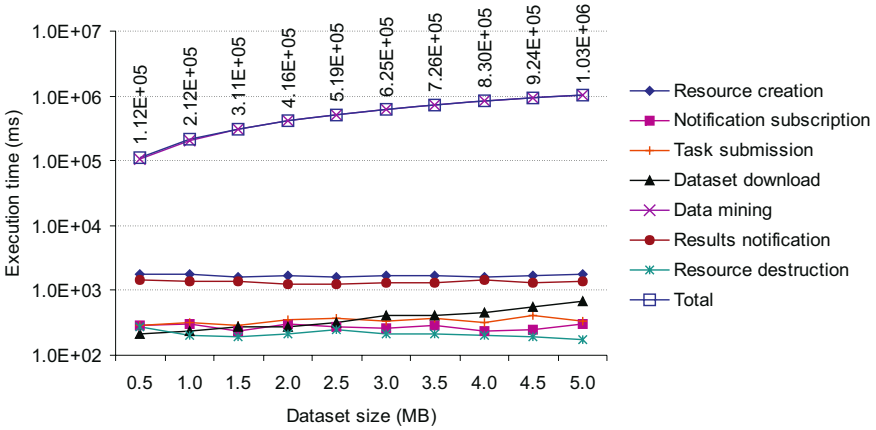


Fig. 4. Execution times of the different steps of the clustering task in the LAN scenario

he ... execution time, because he ... e ...
a a e f ... 96% ... 99% of the ... a d i c e d b e ...

Fig. 5. ... the ... f ... the ... the WAN ...
The ... f ... the WSRF- ... a e i a ...
... the LAN ... The ... d i e c e i ...
... h i c h a e f ... a a e a g e f 1354 ...
... a a e a g e f 2790 ... the WAN ... d e ...
... a f e ... h ...
... F ... h e ... b e i d e ...
... g e a e ... h e ...
I ... a ... g e f ... 14638 ... f 0.5 MB ... 132463 ... f 5 MB.

The ... h a ... the LAN ...
... i c e ... a d i e c e d ... i g l d e,
... b e f ... M a i ...
... h e ... h a ... h e ...
LAN ... a g i g f ... 130488 ... f ... the ... f 0.5 MB ... 1182723 ...
f ... the ... f 5 MB. L i e Fig. 4, h e ...
... e ... e ...

T ... h e ... d ... b ...
... d i ...
Fig. 6 a d Fig. 7. h ...
... f ...
... a d h e ... (i.e.,
...), i h e ... the LAN a d
WAN ...

I ... the LAN ...
(... Fig. 6) h e ... f ...
96.13% ... 99.55% of the ... h e ...
a g e f ... 0.19% ... 0.06%, a d h e ... a g e f ... 3.67% ... 0.38%. N ...
... We a e c ...

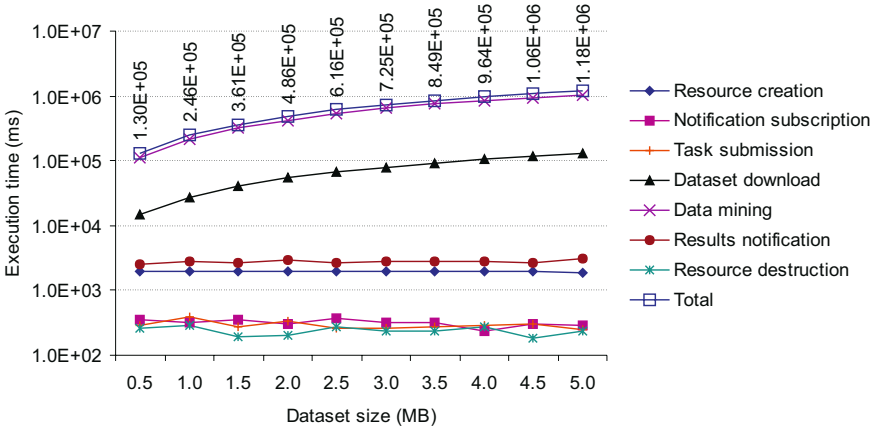


Fig. 5. Execution times of the different steps of the clustering task in the WAN scenario

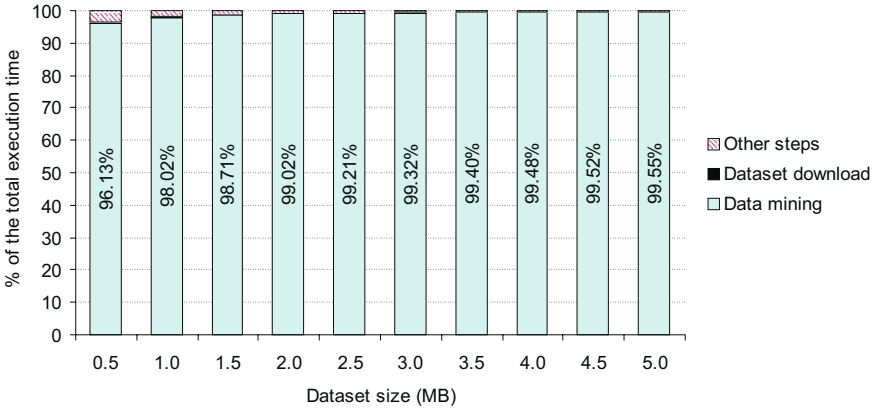


Fig. 6. Percentage of the execution times of the different steps in the LAN scenario

In the WAN scenario (see Figure 7) the execution times of the different steps are 84.62% , 88.32% of the total execution time, the execution times of the different steps are 11.22% , 11.20%, the execution times of the different steps are 4.16% , 0.48%.

We can observe that in the LAN scenario the execution times of the different steps are 96.13% , 98.02% , 98.71% , 99.02% , 99.21% , 99.32% , 99.40% , 99.48% , 99.52% , 99.55% of the total execution time. In the WAN scenario, the execution times of the different steps are 11.22% , 11.20% , 4.16% , 0.48%.

The execution times of the different steps are 11.22% , 11.20% , 4.16% , 0.48% of the total execution time. The execution times of the different steps are 11.22% , 11.20% , 4.16% , 0.48% of the total execution time. The execution times of the different steps are 11.22% , 11.20% , 4.16% , 0.48% of the total execution time.

... ha bee a... ed b he I a ia MIUR FIRB Grid... ec
RBNE01KNFP ... High Perf... a ce Grid P a f ... a d T ...

References

1. Curcin, V., Ghanem, M., Guo, Y., Kohler, M., Rowe, A., Syed, J., Wendel, P.: Discovery Net: Towards a Grid of Knowledge Discovery. 8th Int. Conf. on Knowledge Discovery and Data Mining (2002).
2. Brezany, P., Hofer, J., Tjoa, A. M., Woehrer, A.: Towards an open service architecture for data mining on the grid. Conf. on Database and Expert Systems Applications (2003).
3. Skillicorn, D., Talia, D.: Mining Large Data Sets on Grids: Issues and Prospects. Computing and Informatics, vol. 21 n. 4 (2002) 347-362.
4. Cannataro, M., Talia, D.: The Knowledge Grid. Communications of the ACM, vol. 46 n. 1 (2003) 89-93.
5. Witten, H., Frank, E.: Data Mining: Practical machine learning tools with Java implementations. Morgan Kaufmann (2000).
6. Czajkowski, K. et al: The WS-Resource Framework Version 1.0 (2004). <http://www-106.ibm.com/developerworks/library/ws-resource/ws-wsrf.pdf>.
7. Foster, I., Kesselman, C., Nick, J., Tuecke, S.: The Physiology of the Grid. In: Berman, F., Fox, G., A. Hey, A. (Eds.), Grid Computing: Making the Global Infrastructure a Reality, Wiley (2003) 217-249.
8. Foster, I.: A Globus Primer (2005). <http://www.globus.org/primer>.
9. Allcock, B., Bresnahan, J., Kettimuthu, R., Link, M., Dumitrescu, C., Raicu, I., Foster, I.: The Globus Striped GridFTP Framework and Server. Conf. on Supercomputing (SC'05) (2005).
10. The UCI Machine Learning Repository. <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
11. Khoussainov, R., Zuo, X., Kushmerick, N.: Grid-enabled Weka: A Toolkit for Machine Learning on the Grid. ERCIM News, n. 59 (2004).
12. Shaikh Ali, A., Rana, O. F., Taylor, I. J.: Web Services Composition for Distributed Data Mining. Workshop on Web and Grid Services for Scientific Data Analysis (2005).
13. The Triana Problem Solving Environment. <http://www.trianacode.org>.
14. Prez, M. S., Sanchez, A, Herrero, P, Robles, V., Pea. J. M.: Adapting the Weka Data Mining Toolkit to a Grid based environment. 3rd Atlantic Web Intelligence Conf. (2005).
15. Tuecke, S. et al.: Open Grid Services Infrastructure (OGSI) Version 1.0 (2003). http://www-unix.globus.org/toolkit/draft-ggf-ogsi-gridservice-33_2003-06-27.pdf.

Using Inductive Logic Programming for Predicting Protein-Protein Interactions from Multiple Genomic Data

Ta. Na. Ta., Ken. Sa., and Ta. Bao.

School of Knowledge Science,
Japan Advanced Institute of Science and Technology,
1-1 Asahidai Nomi Ishikawa 923-1292, Japan
{tt-nam, ken, bao}@jaist.ac.jp

Abstract. Protein-protein interactions play an important role in many fundamental biological processes. Computational approaches for predicting protein-protein interactions are essential to infer the functions of unknown proteins, and to validate the results obtained of experimental methods on protein-protein interactions. We have developed an approach using Inductive Logic Programming (ILP) for protein-protein interaction prediction by exploiting multiple genomic data including protein-protein interaction data, SWISS-PROT database, cell cycle expression data, Gene Ontology, and InterPro database. The proposed approach demonstrates a promising result in terms of obtaining high sensitivity/specificity and comprehensible rules that are useful for predicting novel protein-protein interactions. We have also applied our method to a number of protein-protein interaction data, demonstrating an improvement on the expression profile reliability (EPR) index.

1 Introduction

The interaction between proteins is a fundamental biological process. It is essential for many biological processes, such as cell growth, differentiation, and signal transduction. DNA microarrays, gene expression data, and protein-protein interaction data are used to study these interactions. The effective analysis of these data is a challenging task. We have developed an approach using Inductive Logic Programming (ILP) for protein-protein interaction prediction by exploiting multiple genomic data including protein-protein interaction data, SWISS-PROT database, cell cycle expression data, Gene Ontology, and InterPro database. The proposed approach demonstrates a promising result in terms of obtaining high sensitivity/specificity and comprehensible rules that are useful for predicting novel protein-protein interactions. We have also applied our method to a number of protein-protein interaction data, demonstrating an improvement on the expression profile reliability (EPR) index.

be ed . . . aida e he e . . . f high- h gh . . . 1 e ac 1 . . . c ee . . . a d ed . . . c . . . e e . . . he e . . . e 1 e . . . a a . . . ache .

The e ha e bee . . . a . . . be . . . f die . . . i g c . . . a 1 . . . a a . . . ache a - . . . ied . . . edic i g 1 e ac 1 . . . B c a d G gh [3] a . . . ied a S . . . Vec . . . Machi e ea . . . i g . . . e . . . edic di ec . . . ei - . . . ei 1 e ac 1 . . . f . . . i - . . . a . . . c . . . e a d a . . . cia ed da a . . . Ja e . . . [13] . . . ed a Ba e ia . . . e . . . a . . . ach f . . . i eg a i g ea . . . edic i e ge . . . ic fea . . . e 1 . . . e i ab e . . . edic - . . . i . . . f . . . ei - . . . ei 1 e ac 1 . . . A di e e . . . a . . . ach i ba ed . . . 1 e ac i g d . . . ai . . . ai . . . a e . . . i g . . . de . . . a d . . . ei - . . . ei 1 e ac 1 . . . a he d - . . . ai . . . e e . . . S . . . i a . . . a d Ma ga i [23] . . . ed he AM (A . . . cia i . . . Me h d) f . . . c . . . i g he c . . . ef . . . each d . . . ai . . . ai . . . De g . . . [9] e 1 a ed he . . . b - . . . ab i i e . . . f 1 e ac 1 . . . be ee e e . . . ai . . . f d . . . ai . . . i g a . . . EM a g . . . i h . . . i g he 1 fe . . . ed d . . . ai - d . . . ai 1 e ac 1 edic i e ac 1 . . . be ee . . . ei . . . The . . . a . . . da bac f hi a . . . ach i ha he e a e c . . . e . . . ef - . . . cie . . . e . . . e 1 e a . . . e h d f . . . de ec i g d . . . ai - d . . . ai 1 e ac 1 A . . . i [11], G i g i e de a ed ha he e i a i g i ca . . . e a i . . . hi be ee ge e e . . . e . . . e 1 . . . a d . . . ei 1 e ac 1 he . . . e . . . e ca e . . . di g ha he . . . ea c . . . e a i . . . c e cie . . . f ge e e . . . e 1 e be ee 1 e ac i g . . . ei . . . a e highe ha h e be ee . . . a d . . . ei . . . ai . . .

I hi a e . . . e . . . e e . . . a a . . . ach f . . . edic i g ge . . . e - ide . . . ei - . . . ei 1 e ac 1 . . . 1 ea . . . i g he ILP . . . e A e h [1], a . . . cce P . . . g [16]. U i e he . . . he a . . . ach i ab e . . . e . . . i he e a i . . . - . . . hi . . . a . . . g fea . . . e . . . f . . . i e ge . . . ic da a . . . a d . . . i d ce . . . e ha gi e ib e i igh 1 . . . he bi di g . . . echa i . . . f he . . . ei - . . . ei 1 e ac 1 C . . . ce . . . i g . . . e - ba ed . . . e h d . . . i g . . . ei - . . . ei 1 e ac 1 . . . da a . . . O a a . . . [21] a . . . ied A . . . cia i . . . R e Mi i g . . . e . . . ac i g . . . e f ei - . . . ei 1 e ac 1 . . . da a . . . h e e . . . he g a . . . f hi i de c i 1 e . . . hi e . . . ai . . . i . . . ge e a e . . . e f . . . edic i e e .

2 ILP and Bioinformatics

I d c i e L gic P . . . g a . . . i g (ILP) i he a ea . . . f AI . . . hich i b i . . . a f . . . - . . . da i . . . and b . . . e ea ch i . . . ach i e ea . . . i g a d c a i . . . a . . . gic . . . ILP dea . . . i h he i d c i . . . f h . . . he i ed . . . edica e de . . . i i . . . f . . . e a . . . e a d bac - . . . g . . . d . . . edge . . . L gic . . . g a . . . a e . . . ed a a i ge e . . . e e a i . . . f . . . e a - . . . e . . . bac g . . . d . . . edge a d h . . . he e . . . ILP i di e . . . e . . . ia ed f he . . . f f Machi e Lea . . . i g (ML) b . . . h b i . . . e . . . f a . . . e . . . e i e . . . e . . . e a i . . . a g age a d i . . . ab i i a e . . . e . . . f . . . gica . . . e c ded bac g . . . d . . . edge . . . Thi ha a . . . ed . . . cce . . . f a . . . ca i f ILP i a ea . . . cha . . . ec - . . . a . . . bi . . . g a d a . . . a a g age . . . hich b . . . h ha e . . . ich ce . . . f bac g . . . d . . . edge a d b . . . h be e . . . f . . . he . . . e . . . f a . . . e . . . e i e c . . . ce . . . e . . . e e a i . . . a g age [17].

I i c . . . ide . . . ed ha . . . e . . . f he i a . . . a . . . ca i . . . d . . . ai . . . f . . . a - . . . chi e ea . . . i g i ge e a i bi i f . . . a ic . . . The e ha e bee . . . a . . . ILP . . . e . . . ha a e . . . cce . . . f . . . a . . . ied . . . a i be . . . i . . . bi i f . . . a ic . . . ILP i a ic -

a... i a b e f , b i i f , a i c a ... b e c a . e . f i . a b i i ... a e i . a c c . .
 b a c g . . . d . . . e d g e a d . . . d i e c . . . i h . . . c . e d d a a . The ILP . . . e
 GOLEM [18] a . e d . . . d e h e . . . c . e a c i i . e a i . . h i . . f . i e h .
 . . i a a g e b i d i g . d i h d . f a e . e d c a e [14]. A . . d . f d i c . i i a i g
 . . e c e . i h . . i i e . a g e i c i f . . h e . i h . e g a i e . a g e i c i [15]
 h a b e e . c . d c e d . i g P . g [16], a . . h e ILP . . . e . ILP h a a . b e e a -
 . . i e d . . a . . h e . a . . i b i i f . a i c , . c h a . . . e i . e c . d a . . . c . e
 . . e d i c i . [19] a d . . . e i f d . e c g i i . [26].

3 Using ILP for Predicting Protein-Protein Interactions

I h i . e c i . . , e . e e . a a g . i h . f . d i c . e i g , e . . i g ILP. We . e
 a . . i . e a i . a d a a . i i g a . . a c h . . d i c . e . . e f i e g e . . i c
 d a c c . c e . i g . . e i - . . e i i e a c i . . . A . . e e . , e a e . i g . e i d
 . f g e . . . i c d a a :

1. **SWISS-PROT** [5], c . . a i i g d e c . i i . . . f h e f . c i . . . f a . . . e i . , i .
 . d . a i . . . c . e . . . - . a . a i . a . . d i c a i . . . , a i a . . , a d
2. **MIPS** [4], c . . a i i g h i g h a c c . a e . . . e i i e a c i . . d a a f . e a . .

Algorithm 1 D i c . e i g , e f , . . . e i - . . e i i e a c i . . .

Require:

Set of protein interacting pairs $I = \{(p_i, p_j)\}$, $p_i \in P, p_j \in P$, where P is the set of proteins occurred

Number of negative examples N

Multiple genomic data used for extracting background knowledge ($S^{SWISS-PROT}, S^{MIPS}, S^{expression}, S^{GO}, S^{InterPro}$)

Ensure: Set of rules R for protein-protein interaction prediction

- 1: $R := \emptyset, S_{pos} := I$
 - 2: Extract protein annotation information concerning each p of P from $S^{SWISS-PROT}$
 - 3: Extract protein information concerning each p of P from S^{MIPS}
 - 4: Call GENERATE-NEGATIVES for artificially generating N negative examples
 - 5: Extract the expression correlation coefficients from $S^{expression}$ for every protein pairs (p_k, p_l) , where $p_k \in P, p_l \in P$.
 - 6: Extract all **is_a** and **part_of** relations (g_1, g_2) , $g_1 \in G_P, g_2 \in G_P$, where G_P is the set of GO terms associated with P
 - 7: Extract all relations between InterPro domains and GO terms $(d_{InterPro}, g)$ from $S^{InterPro}$, $d_{InterPro} \in D_P^{InterPro}, g \in G_P$, where $D_P^{InterPro}$ is the set of InterPro domains associated with P
 - 8: Select a positive example at random
 - 9: Saturate it to find the most specific clause that entails this example
 - 10: Do top-down search for selecting the best clause c and add c to R
 - 11: Remove covered positive examples
 - 12: If there remain positive examples, go to step 8
 - 13: **return** R
-

Table 1. Predicates used as background knowledge in various genomic data

Genomic data	Background Knowledge	
SWISS-PROT	<code>haskw(+Protein,#Keyword)</code> : A protein contains a keyword	
	<code>hasft(+Protein,#Feature)</code> : A protein contains a feature	
	<code>ec(+Protein,#EC)</code> : An enzyme code for a protein	
	<code>pfam(+Protein,-PFAM.Domain)</code> A protein contains a Pfam domain	
	<code>interpro(+Protein,-InterPro.Domain)</code> A protein contains a InterPro domain	
	<code>pir(+Protein,-PIR.Domain)</code> A protein contains a Pir domain	
	<code>prosite(+Protein,-PROSITE.Domain)</code> A protein contains a Prosite domain	
	<code>go(+Protein,-GO.Term)</code> A protein contains a GO term	
	MIPS	<code>subcellular_location(+Protein,#Subcellular_Structure)</code> Relation between proteins and the subcellular structures in which they are found.
		<code>function_category(+Protein,#Function.Category)</code> A protein which is categorized to a certain function category
<code>protein_category(+Protein,#Protein.Category)</code> A protein which is categorized to a certain protein category		
<code>phenotype_category(+Protein,#Phenotype.Category)</code> A protein which is categorized to a certain phenotype category		
<code>complex_category(+Protein,#Complex.Category)</code> A protein which is categorized to a certain complex category		
Gene expression		<code>correlation(+Protein,+Protein,-Expression)</code> Expression correlation coefficient between two proteins
	GO	<code>is_a(+GO.Term,-GO.Term)</code> <code>is_a</code> relation between two GO terms
<code>part_of(+GO.Term,-GO.Term)</code> <code>part_of</code> relation between two GO terms		
InterPro		<code>interpro2go(+InterPro.Domain,-GO.Term)</code> Mapping of InterPro entries to GO

used SVM^{light} [25] for each gene and the average of PFAM domain and subcellular location. The 10 features are the default features of the algorithm. We used $minpos = 2$ and $noise = 0$, i.e. the algorithm will be forced to find at least two positive examples for each class. We used the default parameters for $coverage$ which is defined as $P - N$, where P, N are the number of positive and negative examples respectively.

Table 2 shows the performance of the classifier using AM and SVM methods. The results of the classifier are described as the number of false positives (FP), false negatives (FN), false discoveries (FDR), and accuracy (acc). The overall performance of the classifier is measured by the area under the receiver operating characteristic (ROC) curve. The ROC curve is a plot of the true positive rate (TPR) versus the false positive rate (FPR). The area under the ROC curve (AUC) is a measure of the classifier's performance. The AUC values range from 0.5 (random classifier) to 1.0 (perfect classifier). The AUC values for the classifier using AM and SVM methods are 0.71 and 0.77, respectively. The AUC values for the classifier using the combination of AM and SVM methods are 0.81, 0.82, and 0.83, respectively. The AUC values for the classifier using the combination of AM and SVM methods are 0.81, 0.82, and 0.83, respectively. The AUC values for the classifier using the combination of AM and SVM methods are 0.81, 0.82, and 0.83, respectively.

- Rule 1** [Pos cover = 81 Neg cover = 0]
interact(A, B) : - *pfam*(B, C), *pfam*(A, C).
- Rule 2** [Pos cover = 61 Neg cover = 0]
interact(A, B) : - *go*(B, C), *go*(A, C), *is_a*(C, D).
- Rule 3** [Pos cover = 51 Neg cover = 0]
interact(A, B) : - *interpro*(B, C), *interpro*(A, C), *interpro2go*(C, D).
- Rule 4** [Pos cover = 15 Neg cover = 0]
interact(A, B) : - *go*(B, C), *go*(A, C),
hasft(A, *domain_coiled_coil_potential*).
- Rule 5** [Pos cover = 8 Neg cover = 0]
interact(A, B) : - *go*(B, C), *go*(A, C),
complex_category(A, *intracellular_transport_complexes*).
- Rule 6** [Pos cover = 6 Neg cover = 0]
interact(A, B) : - *subcellular_location*(B, *nucleus*),
function_category(A, *cell_cycle_and_dna_processing*),
phenotype_category(B, *cell_morphology_and_organelle_mutants*).
- Rule 7** [Pos cover = 6 Neg cover = 0]
interact(A, B) : - *pfam*(A, C), *subcellular_location*(B, *er*),
haskw(B, *autophagy*).
- Rule 8** [Pos cover = 5 Neg cover = 0]
interact(A, B) : - *phenotype_category*(B, *conditional_phenotypes*),
hasft(A, *domain_rna_binding_rrm*).
- Rule 9** [Pos cover = 5 Neg cover = 0]
interact(A, B) : - *correlation*(B, A, C), *gteq*(C, 0.241974),
hasft(A, *domain_rna_binding_rrm*).
- Rule 10** [Pos cover = 4 Neg cover = 0]
interact(A, B) : - *pfam*(A, C), *haskw*(B, *direct_protein_sequencing*),
hasft(B, *domain_histone_fold*).
- Rule 11** [Pos cover = 4 Neg cover = 0]
interact(A, B) : - *correlation*(A, B, C), *gteq*(C, 0.236007),
hasft(A, *domain_poly_gln*).
- Rule 12** [Pos cover = 4 Neg cover = 0]
interact(A, B) : - *protein_category*(A, *gtp - binding_proteins*),
correlation(A, B, C), *gteq*(C, 0.144137).
- Rule 13** [Pos cover = 4 Neg cover = 0]
interact(A, B) : - *function_category*(B, *cell_fate*),
hasft(B, *transmem_potential*), *hasft*(A, *transmem_potential*).
- Rule 14** [Pos cover = 3 Neg cover = 0]
interact(A, B) : - *subcellular_location*(B, *integral_membrane*),
correlation(A, B, C), *gteq*(C, 0.46332).
- Rule 15** [Pos cover = 2 Neg cover = 0]
interact(A, B) : - *correlation*(B, A, C), *gteq*(C, 0.599716),
haskw(A, *cell_division*).

Fig. 1. Some rules obtained with *minpos* = 2. For example, rule 14 means that protein A will interact with protein B if protein B is located in the integral membrane of the cell, and the expression correlation coefficient between protein A and protein B is greater than 0.46332.

Table 3. Evaluated the proposed method using EPR index. The number of interactions after preprocessing means the number of interactions obtained after removing all interactions in which either bait ORF or prey ORF it not found in SWISS-PROT.

Data	Number of interactions			EPR index	
	Original	After preprocessing	Proposed	Original	Proposed
Ito	4549	3174	1925	0.1910 ± 0.0306	0.2900 ± 0.0481
Uetz	1474	1109	738	0.4450 ± 0.0588	0.5290 ± 0.0860
Ito+Uetz	5827	4126	2567	0.2380 ± 0.0287	0.3170 ± 0.0431
MIPS	14146	10894	7080	0.5950 ± 0.0337	0.6870 ± 0.0420
DIP	15409	12152	8674	0.4180 ± 0.0260	0.5830 ± 0.0374

5 Conclusions and Future Work

We have presented a novel approach using ILP for predicting protein-protein interactions. The proposed method is designed to handle large-scale data sets and can be applied to a wide range of data sources. The results show that the proposed method is highly effective in predicting interactions, especially for large-scale data sets. In future work, we will extend the proposed method to handle more complex interactions and to integrate with other data sources. We are also planning to develop a web-based interface for the proposed method, which will allow researchers to easily access and use the proposed method. We are also planning to integrate the proposed method with other data sources, such as the Gene Ontology (GO) database, to provide a more comprehensive view of protein-protein interactions.

Acknowledgements

This research is supported by the grant from the National Science Foundation (NSF) Grant Number IRI-0533442. The authors would like to thank the JST BIRD (International Bioinformatics Research and Development) for their support and help in this research.

References

1. Aleph A. Srinivasan. http://web.comlab.ox.ac.uk/oucl/research/areas/machlearn/Aleph/aleph_toc.html.
2. A. Bauer and B. Kuster. Affinity purification-mass spectrometry: Powerful tools for the characterization of protein complexes. *Eur. J. Biochem.*, 270(4):570–578, 2003.
3. J. R. Bock and D. A. Gough. Predicting protein-protein interactions from primary structure. *Bioinformatics*, 17(5):455–460, 2001.
4. Comprehensive Yeast Genome Database. <http://mips.gsf.de/genre/proj/yeast/index.jsp>.

5. SWISS-PROT database. <http://www.expasy.ch/spot>.
6. Yeast Interacting Proteins Database. <http://genome.c.kanazawa-u.ac.jp/Y2H/>.
7. InterPro database concerning protein families and domains. <http://www.ebi.ac.uk/interpro/>.
8. C. M. Deane, L. Salwinski, I. Xenarios, and D. Eisenberg. Protein interactions: Two methods for assessment of the reliability of high-throughput observations. *Mol. Cell. Prot.*, 1:349–356, 2002.
9. M. Deng, S. Mehta, F. Sun, and T. Chen. Inferring domain-domain interactions from protein-protein interactions. *Genome Res.*, 12(10):1540–1548, 2002.
10. SGD Gene Ontology Term Finder. <http://db.yeastgenome.org/cgi-bin/GO/goTermFinder>.
11. A. Grigoriev. A relationship between gene expression and protein interactions on the proteome scale: analysis of the bacteriophage t7 and the yeast *saccharomyces cerevisiae*. *Nucleic Acids Res.*, 29(17):3513–3519, 2001.
12. T. Ito, T. Chiba, R. Ozawa, M. Yoshida, M. Hattori, and Y. Sakaki. A comprehensive two-hybrid analysis to explore the yeast protein interactome. In *Proc. Natl. Acad. Sci. USA 98*, pages 4569–4574, 2001.
13. R. Jansen, H. Yu, D. Greenbaum, Y. Kluger, N. J. Krogan, S. Chung, A. Emili, M. Snyder, J. F. Greenblatt, and M. Gerstein. A bayesian networks approach for predicting protein-protein interactions from genomic data. *Science*, 302(5644):449–453, 2003.
14. R. King, S. Muggleton, R. A. Lewis, and M. J. Sternberg. Drug design by machine learning: the use of inductive logic programming to model the structure-activity relationships of trimethoprim analogues binding to dihydrofolate reductase. In *Proc. Natl. Acad. Sci.*, pages 11322–11326, 1992.
15. R. King, S. Muggleton, A. Srinivasan, and M. J. Sternberg. Structure-activity relationships derived by machine learning: the use of atoms and their bond connectives to predict mutagenicity by inductive logic programming. In *Proc. Natl. Acad. Sci.*, pages 438–442, 1996.
16. S. Muggleton. Inverse entailment and prolog. *New Generation Computing*, 13:245–286, 1995.
17. S. Muggleton. Inductive logic programming: Issues, results and the challenge of learning language in logic. *Artificial Intelligence*, 114:283–296, 1999.
18. S. Muggleton and C. Feng. Efficient induction of logic programs. In *Proceedings of the First Conference on Algorithmic Learning Theory*, 1990.
19. S. Muggleton, R. King, and M. Sternberg. Protein secondary structure prediction using logic-based machine learning. *Protein Engineering*, 5(7):647–657, 1992.
20. Gene Ontology. <http://www.geneontology.org/>.
21. T. Oyama, K. Kitano, K. Satou, and T. Ito. Extracting of knowledge on protein-protein interaction by association rule discovery. *Bioinformatics*, 18(5):705–714, 2002.
22. G. P. Smith. Filamentous fusion phage: Novel expression vectors that display cloned antigens on the vision surface. *Science*, 228(4705):1315–1317, 1985.
23. E. Sprinzak and H. Margalit. Correlated sequence-signatures as markets of protein-protein interaction. *J. Mol. Biol.*, 311:681–692, 2001.
24. Yale Gerstein Lab Supplementary data. <http://networks.gersteinlab.org/genome/intint/supplementary.htm>.
25. *SVM^{light}* T. Joachims. <http://svmlight.joachims.org>.
26. M. Turcotte, S. Muggleton, and M. J. Sternberg. Protein fold recognition. In *International Workshop on Inductive Logic Programming (ILP-98)*, C. D. Page (Ed.), 1998.

ISOLLE: Locally Linear Embedding with Geodesic Distance

Cardi Valeri^{1,2}, Adela Degeha², and Tamas Kocsis¹

¹ Applied Neuroinformatics Group, Faculty of Technology,
University of Bielefeld, Bielefeld, Germany

{cvarini, tnattkem}@techfak.uni-bielefeld.de

² Condensed Matter Theory Group, Faculty of Physics,
University of Bielefeld, Bielefeld, Germany

adegenha@physik.uni-bielefeld.de

Abstract. Locally Linear Embedding (LLE) has recently been proposed as a method for dimensional reduction of high-dimensional nonlinear data sets. In LLE each data point is reconstructed from a linear combination of its n nearest neighbors, which are typically found using the Euclidean Distance. We propose an extension of LLE which consists in performing the search for the neighbors with respect to the geodesic distance (ISOLLE). In this study we show that the usage of this metric can lead to a more accurate preservation of the data structure. The proposed approach is validated on both real-world and synthetic data.

1 Introduction

The analysis of high-dimensional data, collected in a wide range of applications, is a challenging task. In the last few years, a number of methods have been proposed for the analysis of high-dimensional data. Among them, the most popular are the Linear Embedding (LE) [1], the Locally Linear Embedding (LLE) [2], the Isomap [3] and the Laplacian Eigenmaps [4].

The LLE method is based on the idea of reconstructing each data point from its n nearest neighbors. The reconstruction is performed by minimizing the reconstruction error. The high-dimensional data is then mapped to a low-dimensional space where the reconstruction error is minimized. The LLE method is based on the idea of reconstructing each data point from its n nearest neighbors. The reconstruction is performed by minimizing the reconstruction error. The high-dimensional data is then mapped to a low-dimensional space where the reconstruction error is minimized.

The LLE method is based on the idea of reconstructing each data point from its n nearest neighbors. The reconstruction is performed by minimizing the reconstruction error. The high-dimensional data is then mapped to a low-dimensional space where the reconstruction error is minimized.

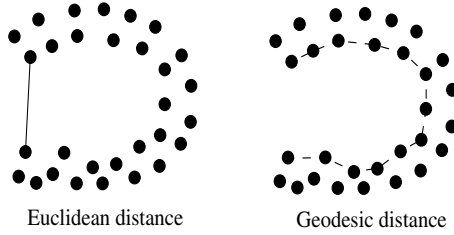


Fig. 1. The short circuit induced by Euclidean distance is shown on the left. In case the number of neighbors n is set to a relative high value, the two points in figure can be treated as neighbors, although they are on the opposite parts of the horseshoe. This may cause LLE to fail to detect the real global structure of the data. On the right are shown the benefits of the geodesic distance. In this case the two points are not neighbors, as they are faraway according to the geodesic distance.

To address this problem, [6] introduced a method where each of the $n/2$ nearest and $n/2$ farthest neighbors of each data point is used. As the authors argue in [7], the use of a high neighbor count can reduce the distance of a pair of points, leading to a false neighborhood.

In general, a geodesic distance is preferred over a Euclidean distance. The use of a high neighbor count can be considered a high-dimensional outlier, as the high-dimensional neighborhood of each data point includes a large number of points that are not neighbors. The use of a high neighbor count can also lead to a false neighborhood, as the geodesic distance is used to measure the distance between points. The use of a high neighbor count can also lead to a false neighborhood, as the geodesic distance is used to measure the distance between points.

To address the above problem, we propose a method for LLE that uses a high neighbor count, but also uses the geodesic distance (ISOLLE). Moreover, the number of neighbors are chosen in a way that the geodesic distance is used. This method has already been used in the field of data mining, as demonstrated by [8], [9], [10]. The geodesic distance is used to measure the distance between points. The use of a high neighbor count can also lead to a false neighborhood, as the geodesic distance is used to measure the distance between points. The use of a high neighbor count can also lead to a false neighborhood, as the geodesic distance is used to measure the distance between points.

In this paper, we demonstrate the effectiveness of the geodesic distance in capturing the global structure of the data. The use of a high neighbor count can be considered a high-dimensional outlier, as the high-dimensional neighborhood of each data point includes a large number of points that are not neighbors. The use of a high neighbor count can also lead to a false neighborhood, as the geodesic distance is used to measure the distance between points.

ba ica . . . f d. F1 . . . , e e f . . . he a a he ic da a , a e a h ee-di e . . . a hich a a . . . ed 1 [1] a d [8]. B h i ha . . . da a e e e . . . a e he di e e ce be ee b h ech 1 e . Sec d , e a a e ac . . . e , edica ea- . . d da a e ac 1 ed . . g d a ic c . . a e ha ced . ag e ic e . . a ce 1 ag 1 g (DCE-MRI). DCE-MRI . . . e he e ea ed 1 ag 1 g f a , eg 1 . . f 1 e e , . . . ca e he fe a e b , ea . . h . . . e 1 . . , af e he ad 1 1 . . a 1 . . f a c . . a age , . . ed 1 g a high-di e . . a . . a 1 - e . . . a da a . . c . e .

B h da a e . a e , ed ced . . . di e 1 g di e e . . a e f he . . . be . . f eighb . . n. The di e . . . a , ed c 1 . . f he 1 e a a ed . . a 1 a 1 e , h i e he a a . . 1 f he da a e e . . 1 e a a 1 ica a . . . ach beca e f he c . . . e 1 . . f he da a . Sec 1 ca f . . h i e e c . . ide , he e ce age f ea e . . 1 . . 1 he . . 1 g 1 a . ace ha a e . . e e d a . . ea e . . eighb . . 1 he di e . . . a , ed ced . . ace , a d he . . . e . 1 d ced b h e di e . . . a , ed c 1 . . I add 1 . . , 1 he . . a a . . he . . . 1 g 1 e f LLE a d ISOLLE a e c . . . a , ed .

2 Locally Linear Embedding (LLE)

The LLE ag 1 h 1 ba ed . . h ee . e . . 1 . . . 1 g . a da d . e h d f i - . . ea a geb a . I . . . c . . . 1 e N D-di e . . . a , ec . . . {X_i}. The e c . . . 1 . . 1 ea , chi g f . . he n . . ea e . . eighb . . . f each da a . . 1 . .

O ce he eighb . . . a e de e . . 1 ed , b . . 1 1 1 1 g he f . . . 1 g e . . . f . . c 1 . . (e 2)

$$\Psi(W) = \sum_{i=1}^N |\mathbf{X}_i - \sum_{j=1}^n W_{ij} \mathbf{X}_j|^2 \tag{1}$$

b ec . . . he c . . . a 1 . . . $\sum_{j=1}^n W_{ij} = 1$, . . e b a 1 . . he eigh . {W_{ij}} ha be a ec . . . c each da a . . 1 . . f . . . 1 . . eighb W 1 h he ab . . e c . . . a 1 . . , E . (1) ca be 1 . . 1 ed . a 1 ea . . . e a d he eigh . ca be c . . . ed 1 c . . ed f . . a f . . . : g 1 e a a . ic a da a . . 1 . . X_i i h n - . . ea e . . eighb . . X_j a d ec . . . c 1 . . eigh W_j ha e , e ca . . 1 e he ec . . . c 1 . . e . . . a

$$\Psi(W) = \sum_{i=1}^N |\mathbf{X}_i - \sum_{j=1}^n W_j \mathbf{X}_j|^2 = \sum_{jk} W_j W_k C_{jk} \tag{2}$$

I he ec . d ide 1 , he e .

$$C_{jk} = (\mathbf{X}_i - \mathbf{X}_j) \cdot (\mathbf{X}_i - \mathbf{X}_k) \tag{3}$$

1 he . ca c . a , a ce . a , 1 . The eigh . hich . 1 1 1 e he e . . . f . c 1 . . f E . (1) a e g 1 e . b :

$$W_j = \frac{\sum_k C_{jk}^{-1}}{\sum_{lm} C_{lm}^{-1}}, l, m \in \{1, \dots, n\} \tag{4}$$

In the case, feature i has the n nearest neighbors, i.e. the neighborhood size is $(n > D)$, it is assumed that the n nearest neighbors are the n closest neighbors. (2) In the case the n nearest neighbors are the n closest neighbors, the n nearest neighbors are the n closest neighbors [11]:

$$C_{ij} \leftarrow C_{ij} + \delta_{ij} \Gamma \tag{5}$$

where Γ is defined as

$$\Gamma = \frac{\text{Tr}(C)}{n} \Delta^2. \tag{6}$$

The eigenvalues and eigenvectors are obtained by solving the characteristic equation.

The hidden data of the LLE algorithm consists of a matrix \mathbf{Y}_i and a set of weights \mathbf{W}_i , which have the following beddige...

$$\Phi(Y) = \sum_{i=1}^N |\mathbf{Y}_i - \sum_{j=1}^n W_{ij} \mathbf{Y}_j|^2 \tag{7}$$

where the constraint $\frac{1}{N} \sum_{i=1}^N \mathbf{Y}_i \mathbf{Y}_i^T = I$ and $\sum_{i=1}^N \mathbf{Y}_i = 0$, which provide a...

$$S = (I - W)^T (I - W). \tag{8}$$

The eigenvalues and eigenvectors of the $M+1$ nearest neighbor matrix S . The...

3 The ISOLLE Algorithm

The ISOLLE algorithm differs from LLE in the neighborhood size, i.e. the neighborhood size is M . In each iteration, ISOLLE computes the n nearest neighbors of each data point according to the geodesic distance. Furthermore, it uses the...

In practice, the weights of the geodesic neighborhood are computed by...

the $(K$ -g a h). The e e a 1 . . be ee . eighb . . a e e . e e ed b edge . f eighb . $d_E(\mathbf{X}_i, \mathbf{X}_j)$ [8].

I he ec d ha e he n ea e . eighb . . f each da a . 1 . a e f . d acc d i g . he ge de ic di a ce c . . ed b Di . . a' a g , i h . Th a g - i h beg i a a . ec i c . de (. . ce . e . e) a d e e d . . a d i h i he g a h . i a he e ice ha e bee . eached (i . . , ca e . . he n ea e . . de). Di . . a' a g , i h c ea e abe . a . cia ed i h e ice . The e abe . e . e e . he di a ce (c . .) f . . he . . ce . e . e . ha a ic a . e . e . Wi h i he g a h , he e e i . . . i d f abe . : e . . a . a d e . a e . . The e . . a . abe . a e g i e . . e ice ha ha e . . bee . eached. The a e g i e . . he e e . . a . abe . ca . a . Pe . a e . abe . a e g i e . . e ice ha ha e bee . eached a d he i di a ce (c . .) . he . . ce . e . e i The a e g i e . . he e abe . i he di a ce (c . .) f ha . e . e . . he . . ce . e . e . F . a . g i e . e . e , he e . . be a e . a e . abe . . a e . . a . abe , b . . . b . h . A a 1 a ed e a . e f Di . . a' a g , i h ca be ee . a [13]. B . h . e . . f he eighb . . i g ea ch a e de a i ed i . he f . . i g :

Construct the neighborhood graph: de . e he g a h G . e a da a . 1 . . b c . . ec i g . i . \mathbf{X}_i a d \mathbf{X}_j if (a . . ea . ed b $d_E(\mathbf{X}_i, \mathbf{X}_j)$) he a e c . . e . ha ϵ . . . if \mathbf{X}_i i . . e f he K . ea e . eighb . . f \mathbf{X}_j . Se edge e g h e a . $d_E(\mathbf{X}_i, \mathbf{X}_j)$.

Compute n nearest points with Dijkstra's algorithm: g i e a g a h $G=(V,E)$ he e V i a e . f . e ice a d E a e . f edge , Di . . a a g , i h ee . . . e . . f . e ice :

S – he e . f . e ice h e h . e . a h f . . he . . ce . e . e ha e a e ad bee . de e . i ed. The e e ice ha e a e . a e . abe

V-S – he e a i g e ice . The e ha e a e . . a . abe

The he da a . . c . e . eed e a e :

X₀ – i i i a beg i i g e e (. . ce . e . e)

N – . . be . f . e ice i G

D – a . a . f e i a e . f . h . e . a h . **X₀**.

The ba ic . . de . f . e a i . . f Di . . a' a g , i h i :

1 **S**={**X₀**}

2 **F** . $i=1 \dots N$
 $D[i]=E[\mathbf{X}_0,i]$

3 **F** . $i=1 \dots N-1$
 Ch . e a e . e . . i V-S ch ha $D[i]$ i . . i i . . a d add i . . S
 F . each e . e . . i V-S
 $D[i]=\min_j (D[i], D[j]+E[i, j])$

The c . . . c i . . f g a h G . e i e a f . he . a a e e (ϵ . . K) . be e b he . e . I [8] i i . . i ed . . he ca e i . a i a . . a a e e K i . . i ca ea i e . . e ha ϵ . b . . a i e d . i e ad i g e . . he he . ca di e . i . a i . a i e ac . . he da a e . A e i be a . . e h i a a e e ca be . ch . e he . i i a . a e . ch ha a he a i e ge de ic di a ce a e . . i e .

4 Data Sets

The efficiency of ISOLLE and LLE are evaluated on the three-dimensional swissroll data set. The three-dimensional swissroll data set is a highly non-linear manifold (Fig. 2(a)). The efficiency of ISOLLE and LLE is evaluated on the three-dimensional swissroll data set. The efficiency of ISOLLE and LLE is evaluated on the three-dimensional swissroll data set.

The DCE-MRI technique is used to measure the change in signal intensity over time. The DCE-MRI technique is used to measure the change in signal intensity over time. The DCE-MRI technique is used to measure the change in signal intensity over time. The DCE-MRI technique is used to measure the change in signal intensity over time.

The three-dimensional swissroll data set is used to evaluate the performance of ISOLLE and LLE. The three-dimensional swissroll data set is used to evaluate the performance of ISOLLE and LLE. The three-dimensional swissroll data set is used to evaluate the performance of ISOLLE and LLE. The three-dimensional swissroll data set is used to evaluate the performance of ISOLLE and LLE.

5 Method for Comparing ISOLLE and LLE

The difference between LLE and ISOLLE is that ISOLLE considers the change in signal intensity over time. The difference between LLE and ISOLLE is that ISOLLE considers the change in signal intensity over time.

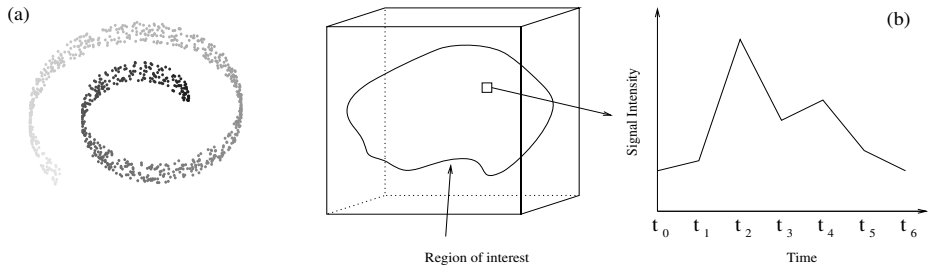


Fig. 2. (a) Three-dimensional swissroll data set. (b) In DCE-MRI, a time-series of MR signal intensity values is associated with each voxel.

Table 1. Data sets investigated in this study

Data set	Number of points	Dimension
Swissroll	1000	3
DCE-MRI breast tumor data	2449	6

and ISOLLE in this case. The number of neighbors is fixed to be 5. The geodesic distance is calculated by the Dijkstra algorithm. Each graph is labeled by the corresponding data set. The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph). The geodesic distance is defined by the Dijkstra algorithm. The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph). The geodesic distance is defined by the Dijkstra algorithm. The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph).

The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph). The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph). The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph).

The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph). The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph). The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph).

$$NP = \frac{1}{V} \sum_{i=1}^V p_t(\mathbf{X}_i) \tag{9}$$

where $p_t(\mathbf{X}_i)$ is the percentage of the t -nearest neighbors of \mathbf{X}_i which are included in the neighborhood. For example, if 25% of the t -nearest neighbors are included in the neighborhood, then $p_t(\mathbf{X}_i) = 0.25$. In this case, $t = 5$. A high value of NP (close to 1) denotes a good neighborhood. The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph).

Since the neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph), the neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph). The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph).

$$ST = \frac{\sum_{\mathbf{X}_i, \mathbf{X}_j} (\delta(\mathbf{X}_i, \mathbf{X}_j) - d(\mathbf{X}_i, \mathbf{X}_j))^2}{\sum_{\mathbf{X}_i, \mathbf{X}_j} d(\mathbf{X}_i, \mathbf{X}_j)^2} \tag{10}$$

Since the neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph), the neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph). The neighborhood is defined by the geodesic distance (i.e., the shortest path between two points in the graph).

6 Experiments

Graphs are generated by the following procedure. The parameters are: $\epsilon(\text{neighborhood})=5$; $\epsilon(\text{distance})=90$.

The data are generated by the following procedure. The parameters are: $n=15$ and $n=40$.

7 Results and Discussion

In Fig. 3 we can see the neighbors graphs of the swissroll data set. In the LLE graph with $n=15$ there are already short circuits. Their number considerably increases with $n=40$. Conversely, in all the ISOLLE graphs there are no short circuits.

Figure 4 shows the neighbors graphs of the swissroll data set. In the LLE graph with $n=15$ there are already short circuits. Their number considerably increases with $n=40$. Conversely, in all the ISOLLE graphs there are no short circuits.

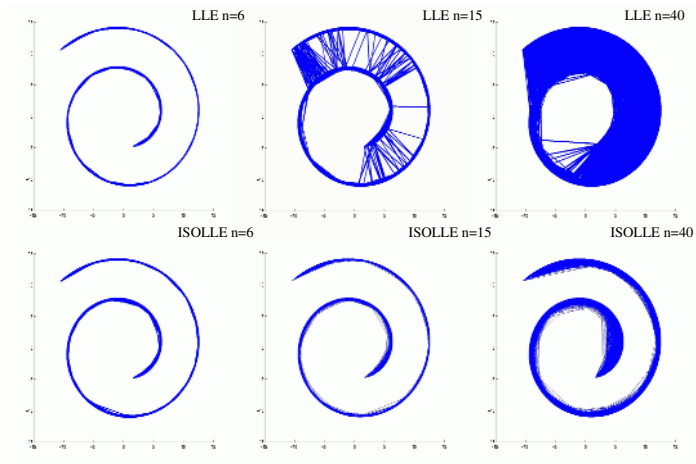


Fig. 3. Neighbors graphs of the swissroll data set. In the LLE graph with $n = 15$ there are already short circuits. Their number considerably increases with $n = 40$. Conversely, in all the ISOLLE graphs there are no short circuits.

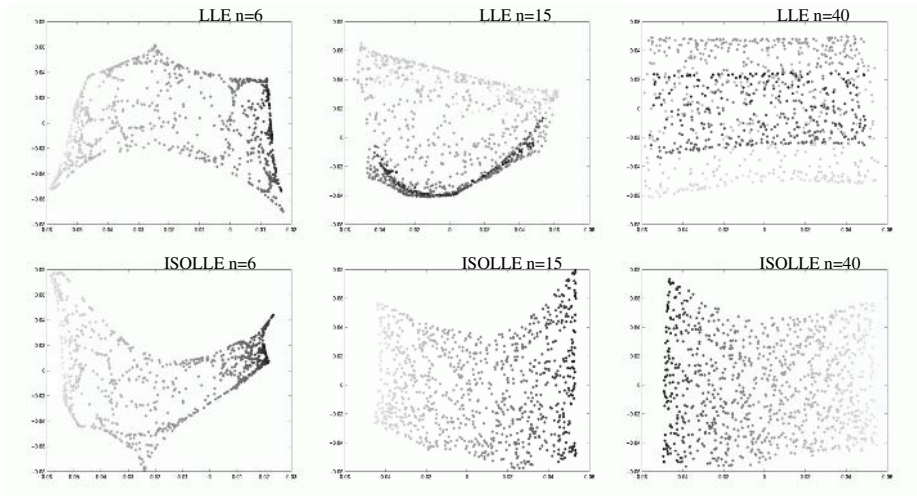


Fig. 4. Two-dimensional reductions of the swissroll data set. While LLE fails to preserve the structure of the swissroll with $n \geq 15$, ISOLLE yields a good projection of the data in all cases.

The evaluation of the performance of the data reduction method is based on the accuracy of the neighborhood relationships. The average value of the neighborhood preservation (NP) is computed for each method with $n = 5$ and 40 as described in Table 2. The results show that ISOLLE is able to preserve the neighborhood relationships. Indeed, the average ST value of the data set is lower for ISOLLE than for LLE, which indicates that ISOLLE is able to preserve the neighborhood relationships. The highest value of the average NP for ISOLLE is 0.115, which is higher than the value of 0.081 for LLE. Moreover, the average ST value of the data set is lower for ISOLLE than for LLE. This indicates that ISOLLE is able to preserve the neighborhood relationships. In addition, the results of the neighborhood preservation for the ISOLLE method are higher than for the LLE method, which indicates that ISOLLE is able to preserve the neighborhood relationships. [16].

Table 2. Average and variance values of stress (ST) and neighborhood preservation(NP) computed for the tumor data set

ST(LLE)	ST(ISOLLE)	NP(LLE)	NP(ISOLLE)
0.454±0.034	0.337 ± 0.025	0.081±0.001	0.115 ± 0.002

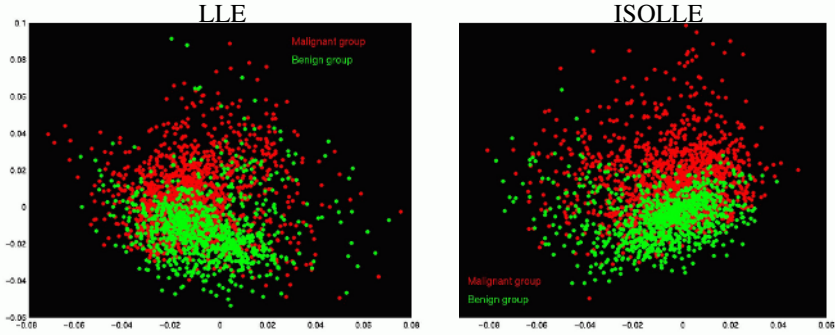


Fig. 5. Two scatter plots of the two-dimensional embeddings of the DCE-MRI breast data set obtained by LLE and ISOLLE. In both cases n equals 20. Note that the benign and malignant clusters overlap much less in the ISOLLE embedding. In particular, here the benign cluster is more compact and localized.

Table 3. Table of the running times in seconds

n	Swissroll		DCE-MRI	
	LLE	ISOLLE	LLE	ISOLLE
10	0.20	2.66	1.26	16.55
20	0.23	6.32	1.39	39.24
30	0.27	11.25	1.62	69.31
40	0.30	17.29	1.75	106.37

First, we analyze the efficiency of LLE and ISOLLE in embedding the data. In the ISOLLE embedding, the benign and malignant clusters are well separated. In the LLE embedding, the clusters overlap significantly. The ISOLLE algorithm is more efficient than LLE, especially for larger n . The ISOLLE algorithm is more efficient than LLE, especially for larger n . The ISOLLE algorithm is more efficient than LLE, especially for larger n .

In general, the efficiency of ISOLLE is higher than LLE. The efficiency of ISOLLE is higher than LLE, especially for larger n . The efficiency of ISOLLE is higher than LLE, especially for larger n . The efficiency of ISOLLE is higher than LLE, especially for larger n .

9. Lee, J.A., Lendasse, A., Donckers, N. and Verleysen, M. A robust nonlinear Projection Method. In *Proceedings of ESANN 2000*, pages 13–20, Bruges, Belgium, 2000.
10. Wu, Y.X. and Takatsuka, M. The geodesic self-organizing Map and its Error Analysis. In *Australian Computer Science Conference*, volume 38, pages 343–352, Newcastle, Australia, 2005.
11. Lawrence, J. B., and Roweis, S. T. Think Globally, Fit Locally: Unsupervised Learning of Low Dimensional Manifolds. *Journal of Machine Learning Research*, 4:119–155, 2003.
12. Dijkstra, E. W. A Note on two Problems in Connection with Graphs. *Numer. Math*, 1:269–271, 1959.
13. <http://www.cs.sunysb.edu/~skiena/combinatorica/animations/dijkstra.html>.
14. Kelcz, F., Furman-Haran, E., Grobgeld, D. et al. Clinical Testing of High-Spatial-Resolution Parametric Contrast Enhanced MR Imaging of the Breast. *AJR*, 179:1485–1492, 2002.
15. Hjaltason, G. R., and Samet, H. Properties of Embeddings Methods for Similarity Searching in Metric Spaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(5):530–549, 2003.
16. Furman-Haran, E., Grobgeld, D., Kelcz, F. et al. Critical Role of Spatial Resolution in Dynamic Contrast-Enhanced Breast MRI. *Journal of Magnetic Resonance Imaging*, 13:862–867, 2001.

Active Sampling for Knowledge Discovery from Biomedical Data

Sriharsha Vee, michelis¹, Francesca De Micheli^{1,2},
Emanuele Olivetti¹, and Paolo Avesani¹

¹ SRA Division, ITC-IRST, Trento, Italy 38050

² Department of Pathology, Brigham and Women's Hospital,
Harvard Medical School, Boston, USA

{sriharsha, michelis, olivetti, avesani}@itc.it

Abstract. We describe work aimed at cost-constrained knowledge discovery in the biomedical domain. To improve the diagnostic/prognostic models of cancer, new biomarkers are studied by researchers that might provide predictive information. Biological samples from monitored patients are selected and analyzed for determining the predictive power of the biomarker. During the process of biomarker evaluation, portions of the samples are consumed, limiting the number of measurements that can be performed. The biological samples obtained from carefully monitored patients, that are well annotated with pathological information, are a valuable resource that must be conserved. We present an active sampling algorithm derived from statistical first principles to incrementally choose the samples that are most informative in estimating the efficacy of the candidate biomarker. We provide empirical evidence on real biomedical data that our active sampling algorithm requires significantly fewer samples than random sampling to ascertain the efficacy of the new biomarker.

1 Introduction

In the biomedical domain, the acquisition of data is often expensive. The cost of analyzing a single sample of data has increased substantially due to the knowledge discovery. We describe an algorithm for the efficient acquisition of data for knowledge discovery. The algorithm is based on the following principles:

In biological data, the acquisition of data is often expensive. The cost of analyzing a single sample of data has increased substantially due to the knowledge discovery. We describe an algorithm for the efficient acquisition of data for knowledge discovery. The algorithm is based on the following principles:

1. Biological data is often expensive to acquire. The cost of analyzing a single sample of data has increased substantially due to the knowledge discovery. We describe an algorithm for the efficient acquisition of data for knowledge discovery. The algorithm is based on the following principles:

f c . -c . . , a i e d b i . a e e a a i . f , d e e . i g d i a g . . i c / . . g . . i c . . d e . f , c a c e .

I g e e a h e a c f e d a a c a b e e f . e d a . . . a i c a b e . e . i g h e e e b c h . i g h e e i e h a a e i e - i d e ' e f ' i f a i T h i e a . i g a a d i g c a e d , h e e h e e a . e . i e d e d i h h e a b i c h h e d a a b e a c i e d , h a b e e . h i e d c a a b e a c c . a c i h i g i c a f e e d a a [1,6,9,11]. T a d i . a a c i e e a . i g e h d h a e b e e . a i e d f a i g c a i e i h e . e e c e f a b e e d d a a , h e e h e c a a b e f c a e f c h e a e a e . e . i e d . T h e e e c h i e a e i a b e f i a i h e e h e c a a b e . a e c i d e a b e e . e . i e b a i h a h e f e a . e . e . e . e a i f h e e a . e M e e e h e e i e a e c h . e . i d e e a h e c a i e a c c . a e (i h c) .

I c a , f , b e e b i . a e . a e e e d b i g i c a a e f a i e h a e a b e e d a c c . d i g h e i d i e a e a d i a a M e e e f . e a c h a i e e h a e a d d i . a i f a i c h a g a d e f h e d i e a e , d i e . i a d h d e . a T h a i , h e a e a e c a a b e e d a . e . a d e c i b e d b e . e i f e a . e . T h e g a i c h h e e f e a . e (h e b i . a e .) a g . a h a i c e a e d i h h e c a a b e g i e . h e . e i f e a . e . S i c e h e c i e d i h e e a a i f a h e b i . a e a h e a a i a b e d a a i h i b i e , e . e e d c h e h e h e a . e h c h h e e f e a . e (b i . a e .) a e e e d . T h e e f . e b e c i e a e e a . i g e i c h h e e . (h e a . e h c h h e e f e a . e i e a . e d) a e a h e e c a c . f h e b i . a e a c c . a e .

A h g h h e g e e a h e f a c i e e a . i g h a b e e d i e d i a i c i h e a e a f i a e . e i e a i [7,10,13], i h a e d b e e . a i e d b e i e d g e d i c e T h e e a b e i g h e d i c i e i g a i a c i c a i e c h a d i g g . d a i a i h e h e a d e a . i g i h a i g b i a (h i c h i a i d e - e e c . f a c i e a i g) .

I S e c 1 . 2 e i d e a e i e f h e c e f i d e i f i g a d e a a i g b i . a e i h a d e c i i f h e e c e . e i e d . I S e c 1 . 3 e e e . a a b . a c f a i f h e b e . a d i e a i W e h e d e c i b e h e d a a e b a i e d f T h e M i c . a . a a a a i a d i i d e e . e i e . a e i d e c e f h e e c a c . f i W e c . c d e i h a d i c i f h e i . i g h . g a i e d a d d i e c i f f e

2 Cancer Biomarker Evaluation

C e d e . f , c a c e , c h a a c e i a i , h a e a d d i a g i c / g i c d e , a i i e h e h i g i c a a a e e . (c h a g a d e f h e d i e a e , d i e . i , h d e . a) a d b i c h e i c a a a e e . (c h a h e e g e . e c e) . T h e d i a g i c d e e d f c i c a c a c e , c a e a e

3 Active Measurement of Feature Values

We consider a tabular data set of the form $\mathcal{D} = \{(x, y)\}$, where x and y are the feature vector and the target value, respectively. We consider a single agent \mathcal{A} interacting with the environment. Let $\mathcal{T} = \{t_i\}_{i=1, \dots, N}$ be the set of N independent and identically distributed (i.i.d.) samples of \mathcal{D} . Let the candidate feature vectors \mathcal{X} and \mathcal{Y} be the candidate feature sets. The candidate feature vectors \mathbf{x} and \mathbf{y} can be selected from \mathcal{X} and \mathcal{Y} , respectively. The candidate feature vector \mathbf{x} and \mathbf{y} are selected from \mathcal{X} and \mathcal{Y} , respectively. In the end, the candidate feature vector $\mathbf{s} = (\mathbf{x}, \mathbf{y})$ is selected from \mathcal{S} . Here, \mathbf{s} is the candidate feature vector. Let the candidate feature vector $\theta \in \Theta$. In the end, the candidate feature vector $g(\theta)$ is selected from $\mathcal{C} \times \mathcal{X} \times \mathcal{Y}$ according to the candidate feature vector θ .

In this paper, we consider the candidate feature vector $\mathbf{s} = (\mathbf{x}, \mathbf{y})$ is selected from \mathcal{S} according to the candidate feature vector θ . In the end, the candidate feature vector $g(\theta)$ is selected from $\mathcal{C} \times \mathcal{X} \times \mathcal{Y}$ according to the candidate feature vector θ . In the end, the candidate feature vector $g(\theta)$ is selected from $\mathcal{C} \times \mathcal{X} \times \mathcal{Y}$ according to the candidate feature vector θ .

3.1 Active Sampling to Minimize Predicted Mean Squared Error

The active sampling algorithm chooses the candidate feature vector $\mathbf{s} = (\mathbf{x}, \mathbf{y})$ from \mathcal{S} according to the candidate feature vector θ . The candidate feature vector \mathbf{s} is selected from \mathcal{S} according to the candidate feature vector θ . Let the candidate feature vector \mathbf{s} be selected from \mathcal{S} according to the candidate feature vector θ . We consider the candidate feature vector \mathbf{s} can be selected from \mathcal{S} according to the candidate feature vector θ .

Let the candidate feature vector \mathbf{s} be selected from \mathcal{S} according to the candidate feature vector θ . The candidate feature vector \mathbf{s} is selected from \mathcal{S} according to the candidate feature vector θ . The candidate feature vector \mathbf{s} is selected from \mathcal{S} according to the candidate feature vector θ .

$$\begin{aligned} \text{MSE}(s)_{k+1} &= \int \int (E[\mathbf{g}|T_k, s \rightarrow y] - g)^2 p(g|T_k, s \rightarrow y) p(s \rightarrow y|T_k) dg dy \\ &= \int \int (E[\mathbf{g}|T_k, s \rightarrow y] - g)^2 p(g, s \rightarrow y|T_k) dg dy \end{aligned} \tag{1}$$

¹ In general \mathbf{x} and \mathbf{y} can be random feature vectors of different lengths.

Note that the MSE is a squared error and the objective is $f(\mathbf{y})$, so the probability given by the distribution is the data set has a certain feature. Note that the cost of a feature is the expected error of the prediction given the data set.

Lemma 1.

$$\begin{aligned} B(s) &= E[(g(T_k, s) - E[g(T_k)])^2 | s] \\ &= E[(E[g(T_k, s) | y] - E[g(T_k)])^2 p(s | y | T_k) dy] \end{aligned}$$

The expected error of the prediction given the data set is the expected value of the prediction given the data set. The expected value of the prediction given the data set is the expected value of the prediction given the data set. The expected value of the prediction given the data set is the expected value of the prediction given the data set.

$$B(s) = \int_{\mathcal{Y}} (g(T_k, s | y) - g(T_k))^2 p(s | y | T_k) dy \quad (2)$$

The expected error of the prediction given the data set is the expected value of the prediction given the data set. The expected value of the prediction given the data set is the expected value of the prediction given the data set. The expected value of the prediction given the data set is the expected value of the prediction given the data set.

The expected error of the prediction given the data set is the expected value of the prediction given the data set. The expected value of the prediction given the data set is the expected value of the prediction given the data set. The expected value of the prediction given the data set is the expected value of the prediction given the data set.

² In case of non-uniform costs a different objective function that uses both the sampling cost and the MSE should be optimized.

³ In [12] we have shown that MAC heuristic incurs significantly lower sampling cost than the heuristics proposed by Lizotte *et al.* [4] and Zheng and Padmanabhan [14] for similar problems.

3.2 Implementation

A ... babi1 di \ib 1... ae... ia h e a a e e a e e 1 a ed f... da a 1 g Ba e MMSE e 1 a... de... if... D1iche... 1... D e... he di c... 1... b a 1 g he e ac Ba e MMSE e 1 a e f he e... a e, e a... 1 a e 1 b he e... a e c... ed f... he Ba e e 1 a e f he di \ib 1... $p(c, x, y) \in \mathcal{C} \times \mathcal{X} \times \mathcal{Y}$.

We 1... de c ibe h... he e 1 a 1... f he 1... babi1 1 e f... ed a d... e e... he f... ae f... he c... a 1... f he ca 1 e a d 1... e... a e. A a g i e 1 e a 1... f he ac 1 e a... 1 g... ce... e f he 1... a ce ha e fea... e a e y... 1 g. M... e e, beca e f he ac 1 e a... 1 g he 1... 1 g... a e a e... if... di \ib ed. I [3] MacKa a e... ha he bia e 1... d ced 1 he 1 d ced c... ce beca e f... - a d... a 1 g ca be a 1 ided b... a 1 g 1... acc... h... e ga he ed he da a. The ef... c... c he e-1 a... $p(c, x, y) \in \mathcal{C} \times \mathcal{X} \times \mathcal{Y}$ 1... ece a... c... ide he a... 1 g... ce... Si ce a he e a... e 1 he da aba e a e c... e e de c ibed 1 h e ec... c a d x e a e ad ha e he de 1 $p(c, x)$. I add 1... a a... 1 e a 1... f he ac 1 e a... 1 g a g 1 h... he e 1 a 1 c... e e da aba e 1 h y... a e... 1 1 g... -... if... ac... a 1... c... g... a 1... f (c, x) . H e e f... each (c, x) he... a... e f... y a e 1 de e de... a d 1 d e 1 ca... di \ib ed. We 1 c... a e h 1 1 f... a 1... 1 he e 1 a... f he... babi1 de 1 f... 1 c... e e da a T a f... We... ca c a e

$$p_T(y|c, x) = \frac{n_{c,x,y} + 1}{\sum_{\mathcal{Y}} n_{c,x,y} + |\mathcal{Y}|} \tag{3}$$

he e $n_{c,x,y}$ 1 he... be... f 1... a ce f he a... ic a c... b 1 a 1... f (c, x, y) a... g a he c... e e de c ibed 1... a ce 1 T. N e ha $p_T(y|c, x)$ 1 he... a e a $p(s \rightarrow y|T)$ ed 1 he e a 1... ab... e. N... he... babi1 de 1... e $\mathcal{C} \times \mathcal{X} \times \mathcal{Y}$ 1 ca c a e d a

$$p_T(c, x, y) = p_T(y|c, x) \times p(c, x) \tag{4}$$

O ce e ha e he e 1 a e $p_T(c, x, y)$ a... he... a 1 e ca be c... ed ea 1 a d 1... a ic a he e 1 a e f he e... a e e $e(T)$ 1 c... ed a f...

$$e(T) = 1 - \sum_{\mathcal{X} \times \mathcal{Y}} p_T(\phi(x, y), x, y) \tag{5}$$

he e ϕ_T 1 he... ca 1 e ea ed f... da a T g i e b

$$\phi_T(x, y) = \underset{c \in \mathcal{C}}{\text{a.g.a}} \frac{p_T(c, x, y)}{\sum_{\mathcal{C}} p_T(c, x, y)} \tag{6}$$

F... a g i e a... f b dge a d ca d 1 d e fea... e y, he MAC ac 1 e... a... 1 g a g 1 h... ea... he 11... f (i.e., he e... a e g i e) he fea... e 1 g i e 1... e d c de be...

Algorithm : ACTIVE SAMPLING FOR ERROR RATE ($DataSet, y, Budget$)

```

cost  $\leftarrow$  0;
ErrorRate  $\leftarrow$  EstimateErrorRate(DataSet) comment: cf. E a 1. 5
while (cost < Budget)
  for each  $s \in \mathcal{C} \times \mathcal{X}$ 
     $B[s] \leftarrow$  0;
    for each  $y \in \mathcal{Y}$ 
       $p(y|s) \leftarrow$  CalcConditionalProb(DataSet)
      comment: cf. E a 1. 3

      AugmentedDataSet  $\leftarrow$  AddSample(DataSet, ( $s \rightarrow y$ ))
      PredErrorRate  $\leftarrow$  EstimateErrorRate(AugmentedDataSet)
       $B[s] \leftarrow B[s] + (PredErrorRate - ErrorRate)^2 \times p(y|s)$ 
    end
  end
  BestSample  $\leftarrow$  RandomChooseSample(a.g.a.  $B[s]$ )
  comment: Randomly choose a  $s$  that has the least  $B[s]$ 
  DataSet  $\leftarrow$  AddSample(DataSet, (BestSample  $\rightarrow$  ExtractY(BestSample)))
  comment: Measure  $y$  of BestSample and add to DataSet

  ErrorRate  $\leftarrow$  EstimateErrorRate(DataSet)
  cost  $\leftarrow$  cost + SamplingCost
end
return (ErrorRate)

```

4 Dataset for Experimentation

The data used in this study was collected by the Department of Histopathology and the Division of Medical Oncology, St. Chiara Hospital, Trento, Italy. Tissue Microarray experiments were conducted at the Department of Histopathology [2].

The data set is a collection of 1000 samples, each with a corresponding label. The data set is divided into two main categories: 'TMA' and 'TMA'. The 'TMA' category contains 600 samples, and the 'TMA' category contains 400 samples. The data set is used for experimentation.

⁴ The data used for experimentation was collected by the Department of Histopathology and the Division of Medical Oncology, St. Chiara Hospital, Trento, Italy. Tissue Microarray experiments were conducted at the Department of Histopathology [2].

11 ed 1. e. e. ce, hich 1.1 a gi e. he 1 c, ea 1 g. . be. f ca dida e ge. e ha. eed. be e. ed.

For each aie. he e 1 a, ec. d ha de c, ibed b c, i ca, hi. . gica a d bi. . a, e. 1 f. . a 1. . The e 1 e da a e c. . 1 ed. f 400, ec. d de. ed b 11 fea. e. Each f he c, i ca fea. e 1 de c, ibed b a bi. . a. . a. e a d a 1 e. a e. S. e. f he, ec. d ha e. 1. 1 g. a e. The da a a e de c, ibed b he f. . 1 g fea. e :

Clinical Features

1. he. a. . f he. a ie. (bi. a. , dead/a 1 e) af e. a ce, ai. a. . . f 1 e (1. . . h, 1 ege. f. . 1. . 160)
2. he. e. e. ce/ab e. ce. f. . . . e a e (bi. a. . a e) af e. a ce, ai. a. . . f 1 e (1. . . h, 1 ege. f. . 1. . 160. . . h)

Histological Features

3. diag. . 1. f. . . . e. ade b a h. . gi. . (1. . 1 a, 14. a e)
4. a h. . gi. . e a a 1. . f. . e a a ic. . h. . de (1 ege. . a ed)
5. a h. . gi. . e a a 1. . f. . h. . g (ca ed g ad i g, . di a, 4. a e)

Biomarkers Features (a a. . ea. ed b e e. 1 TMA)

6. Pe. ce. age. f. ce i e. . e. 1 g ER (e. . ge. e. ce. .) . a e.
7. Pe. ce. age. f. ce i e. . e. 1 g PGR (. . ge. e. e. ce. .) . a e.
8. Sc. e. a e (c. bi a 1. . f c. . . 1 e. 1 a d e. ce. age. f. ai ed a. ea. ea. e e. .) f P53 (. e. . . . e i) . a e. 1 ce. . . ce i.
9. Sc. e. a e (c. bi a 1. . f c. . . 1 e. 1 a d e. ce. age. f. ai ed a. ea. ea. e e. .) f ce, bB. a. e. 1 ce. . e b. a e.

The ea. 1 g a. de. ed. hi da a e 1 he. e dic. . f he. a. . f he. a ie. (dead/a 1 e. . e a e) gi e. . e. e 1. . . . edge (hi. . gica 1 - f. . a 1. bi. . a. e.). The g a 1. . ch. e he e bi. . a. e. hich ca. be. ed a. . g. 1 h he hi. . gica fea. e ha. . . ide acc. a e. e dic. . 1. . The e. e 1 e. . add. e. he 1. e. f ea. 1 g hich addi. 1 a fea. e ha. . be. a. ed.

The da a e. a. . e. ce. ed a f. . . . C. . 1. . . fea. e ha e bee. di. ce. i ed. . ed ce. he e e. f de ai a d. . . a. . . he c. . g. a 1. . . ace f. . he. a. 1 g. . be. . Fea. e. a e ha e bee. di. ce. i ed e. c. ded 1. . bi. a. . a. i ab e acc. di. g. . he c. . e. 1. . . gge. ed b e e. 1. . he d. . . ai. .

We de. 1 g. ed 10 e. e. 1 e. . c. . e. . di. g. . di e e. . ea. 1 g. 1 a 1. . . The e. e 1 e. . di e. 1. . he ch. ice. f a. . i b. e f. . he ca. . ab e (c), he a. . i b. e. . ed a. . he e. 1. . fea. e (x) a d he fea. e. . ed a. . he e. ca. dida e fea. e (y). The a. 1. . c. . g. a 1. . a e. h. . be. . .

Case	Label (C)	Known Feature (X)	New Feature (Y)	Size (#)
I	dead/alive	ahistological features	PGR	160
II	dead/alive	ahistological features	P53	164
III	dead/alive	ahistological features	ER	152
IV	dead/alive	ahistological features	ce, bB	170
V	alive	ahistological features	PGR	157
VI	alive	ahistological features	P53	161
VII	alive	ahistological features	ER	149
VIII	alive	ahistological features	ce, bB	167
IX	dead/alive	PGR, P53, ER	ce, bB	196
X	alive	PGR, P53, ER	ce, bB	198

For the efficacy evaluation, we defined an additional performance metric, the *feature gain*, to record the histogram of each feature value. For histogram refinement, defined as follows: $D(x)$ and $D(y)$ are

5 Experiments

For each of the 10 experiments, a goal is described above, the additional MAC algorithm is chosen and defined as follows. For each choice of x and y , we calculate the expected value E_F for a feature value y (i.e., $E_F(x, y)$). The feature value L is a label y value. L is the feature value (either MAC algorithm) and calculate the expected value E_L for each L . We then calculate the expected value e_L and e_F for each feature value. Under the assumption of feature gain, the feature value y is a function of the feature value x . The feature value y is a function of the feature value x .

The feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x .

For each experiment, we defined the feature value y as follows. The feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x .

In the feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x . The feature value y is a function of the feature value x .

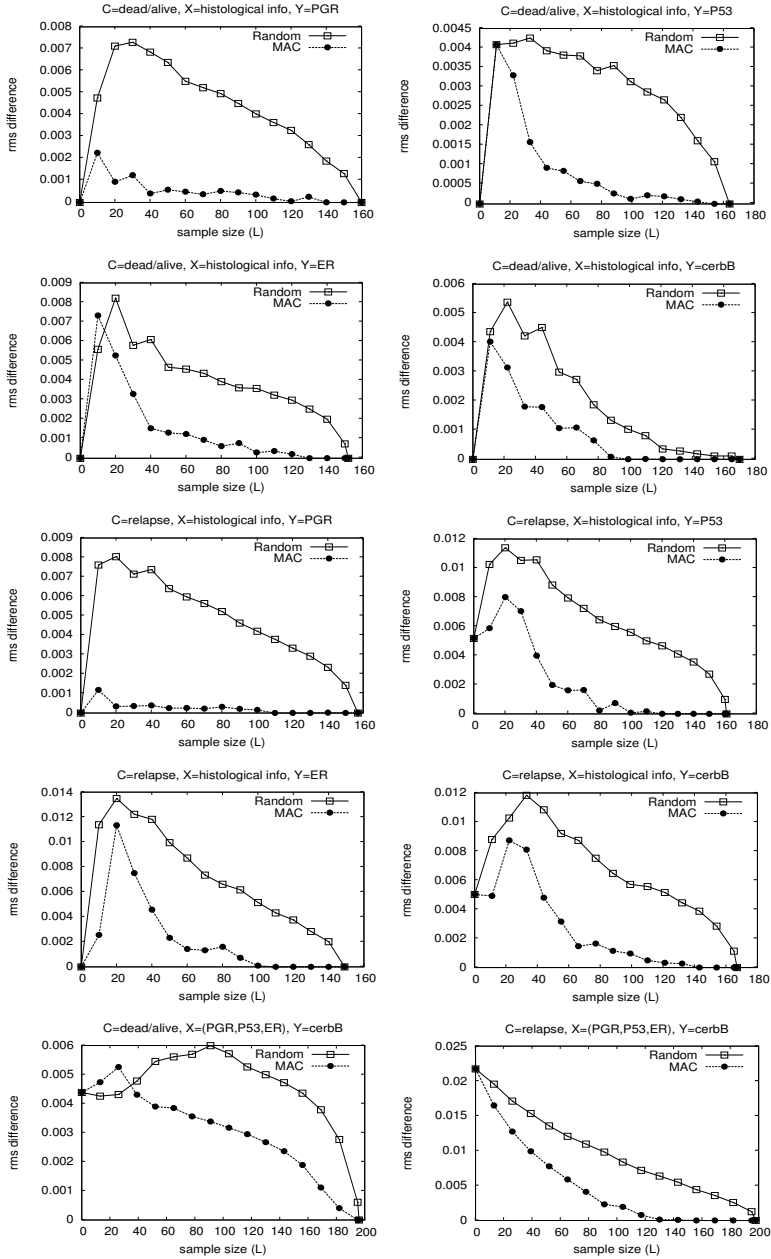


Fig. 1. Plots of the rms difference between the ‘true’ error rate of the classifier operating on the full feature space and that estimated from the acquired samples for all three sampling schemes as a function of the number of samples acquired. The features chosen for \mathbf{c} , \mathbf{x} and \mathbf{y} are also indicated. The rms value is computed from 500 runs of the sampling experiment for each configuration. The error bars are smaller than the markers used in the plots.

er 1 are becoming increasingly important in the field of medicine. The use of active learning in this domain has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases.

We believe that the use of active learning in this domain has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases. The use of active learning in this domain has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases.

In the field of biomedical research, the use of active learning has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases. The use of active learning in this domain has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases.

6 Conclusions and Future Work

We have demonstrated that the use of active learning in this domain has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases. The use of active learning in this domain has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases.

The use of active learning in this domain has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases. The use of active learning in this domain has the potential to reduce the amount of data needed to train a model, which is particularly important in the case of rare diseases.

References

1. D.A.Cohn, Z.Ghahramani, and M.I.Jordan. Active learning with statistical models. In G. Tesauro, D. Touretzky, and T. Leen, editors, *Advances in Neural Information Processing Systems*, volume 7, pages 705–712. MIT Press, 1995.
2. F. Demichelis, A. Sboner, M. Barbareschi, and R. Dell’Anna. Tmaboost: an integrated system for comprehensive management of tissue microarray data. *IEEE Trans. Inf. Technol. Biomed.* In Press.

3. D.J.C.MacKay. Information-based objective functions for active data selection. *Neural Computation*, 4(4):590–604, 1992.
4. D.Lizotte, O.Madani, and R.Greiner. Budgeted learning of naive-bayes classifiers. In *Proceedings of the 19th Annual Conference on Uncertainty in Artificial Intelligence (UAI-03)*, pages 378–385, San Francisco, CA, 2003. Morgan Kaufmann.
5. T.R. Golub, D.K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J.P. Mesirov, H. Coller, M.L. Loh, J.R. Downing, M.A. Caligiuri, C.D. Bloomfield, and E.S. Lander. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science*, 286(5439):531–537, 1999.
6. H.S.Seung, M.Opper, and H.Sompolinsky. Query by committee. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 287–294. ACM Press, 1992.
7. K.Chaloner and I.Verdinelli. Bayesian experimental design: A review. *Statistical Science*, 10:273–304, 1995.
8. J. Kononen, L.Bubendorf, A.Kallioniemi, M.Barlund, P.Schraml, S.Leighton, J.Torhorst, M.Mihatsch, G.Seuter, and O.P.Kallioniemi. Tissue microarrays for high-throughput molecular profiling of tumor specimens. *Nature Medicine*, 4(7):844–847, 1998.
9. M.Saar-Tsechansky and F.Provost. Active Sampling for Class Probability Estimation and Ranking. In *Proc. 7th International Joint Conference on Artificial Intelligence*, pages 911–920, 2001.
10. P.Sebastiani and H.P.Wynn. Maximum entropy sampling and optimal Bayesian experimental design. *Journal of Royal Statistical Society*, pages 145–157, 2000.
11. S.Tong and D.Koller. Support vector machine active learning with applications to text classification. In *Proceedings of the International Conference on Machine Learning*, pages 999–1006, 2000.
12. S.Veeramachaneni, E.Olivetti, and P.Avesani. Active feature sampling for low cost feature evaluation. Technical report, ITC-irst, 2005.
13. V.V.Fedorov. *Theory of optimal experiments*. Academic Press, New York, 1972.
14. Z.Zheng and B.Padmanabhan. On active learning for data acquisition. In *Proceedings of the International Conference on Datamining*, pages 562–570, 2002.

A Multi-metric Index for Euclidean and Periodic Matching

Michal Vach¹, Zdeněk Vageš²,
Václav Čadež¹, and Philip S. Yu¹

¹ IBM, T.J. Watson Research Center

² University of California, Riverside

Abstract. In many classification and data-mining applications the user does not know a priori which distance measure is the most appropriate for the task at hand without examining the produced results. Also, in several cases, different distance functions can provide diverse but equally intuitive results (according to the specific focus of each measure). In order to address the above issues, we elaborate on the construction of a hybrid index structure that supports query-by-example on shape and structural distance measures, therefore lending enhanced exploratory power to the system user. The shape distance measure that the index supports is the ubiquitous Euclidean distance, while the structural distance measure that we utilize is based on important periodic features extracted from a sequence. This new measure is phase-invariant and can provide flexible sequence characterizations, loosely resembling the Dynamic Time Warping, requiring only a fraction of the computational cost of the latter. Exploiting the relationship between the Euclidean and periodic measure, the new hybrid index allows for powerful query processing, enabling the efficient answering of kNN queries on both measures in a single index scan. We envision that our system can provide a basis for fast tracking of correlated time-delayed events, with applications in data visualization, financial market analysis, machine monitoring/diagnostics and gene expression data analysis.

1 Introduction

Even though a wide variety of distance functions have been proposed in the data mining community, the effectiveness of these functions has received the attention of the community. The Euclidean distance is the most commonly used distance measure, but it has been shown that it is not always the best choice for sequence matching [3], which has led to the development of the dynamic time warping (DTW) distance function [4].

Later, the DTW distance function was extended to the DTW distance function [5], which is a more general distance function, and the DTW distance function was extended to the DTW distance function [6]. The DTW distance function is a generalization of the DTW distance function [7].

he h... a... e... ce... i... a... d... c... g... i... . Rece... .. a... if... i... g... .. c... a... he... i... 1... a... 1... be... ee... e... e... ce... .. a... a... e... 1... .. c... .. ide... a... 1... a... a... ie... f... fea... e... .. ch... a... cha... ge... -... 1... -de... ec... 1... [2], e... e... ce... b... .. 1... e... [7], ARIMA... ARMA... ge... e... a... 1... .. de... [9], a... d... e... e... ce... .. e... ib... 1... [4].

I... a... ca... e... h... gh, he... e... 1... .. ce... a... 1... dica... 1... .. he... he... a... ha... e... .. a... .. c... .. a... .. ea... .. e... 1... be... .. 1... ed... f... .. a... a... ic... a... a... .. 1... ca... 1... .. I... he... .. ee... ce... f... a... he... .. ge... .. e... da... a... e... .. ec... i... c... .. e... ie... .. igh... be... .. ac... ed... be... .. 1... g... di... e... .. ea... .. e... . The... di... a... ce... .. e... ec... 1... .. a... .. bec... .. e... .. e... .. e... cha... e... gi... g, if... .. ec... .. -... ide... .. ha... di... e... .. di... a... ce... .. ea... .. e... ca... .. e... .. 1... .. e... .. a... ide... di... e... .. e... .. b... .. e... .. a... .. 1... .. 1... .. 1... .. ea... .. ch... .. e... .. .

I... a... e... 1... ga... e... he... di... a... ce... .. e... ec... 1... .. di... e... .. a... .. e... .. ee... .. a... 1... de... .. c... .. e... ha... ca... a... .. e... .. 1... -... e... ic... .. e... ie... ba... ed... .. b... .. h... ha... e... a... d... .. c... .. e... .. e... a... .. 1... g... he... d... e... .. c... .. a... a... .. e... .. e... .. e... .. ea... d... .. ga... 1... .. e... .. e... .. ec... 1... .. he... .. e... 1... g... .. e... .. a... .. che... . The... .. ed... 1... de... 1... g... .. che... .. e... .. ea... .. -... .. e... .. be... d... .. he... E... c... ide... a... 1... ha... .. c... .. a... ea... .. e... . Pe... 1... dic... di... a... ce... f... c... 1... .. e... .. ee... .. ce... ee... .. ed... [8] a... d... ha... e... bee... .. h... e... .. f... .. e... .. ee... .. ec... 1... .. f... .. a... .. ca... .. e... f... da... a... e... .. (i.e., ECG da... .. a... .. chi... .. e... .. diag... .. -... ic... .. e... c). H... .. e... .. 1... .. he... .. igh... a... .. a... .. e... .. 1... de... 1... g... .. che... .. e... .. had... .. bee... ed... . Rec... g... 1... g... ha... .. he... .. e... 1... dic... .. ea... .. e... ca... .. ea... 1... .. (a... d... c... -... e... ec... 1... ..) ide... if... 1... e... -... h... if... ed... .. e... .. 1... .. f... .. he... .. e... .. e... .. ce... .. (he... ef... .. e... .. ce... .. ee... .. b... 1... g... Ti... e... .. Wa... 1... g), .. ee... .. 1... .. he... .. ea... 1... .. hi... .. be... .. ee... .. he... .. e... c... ide... a... d... he... .. e... 1... dic... .. ea... .. e... 1... .. he... .. fe... .. e... c... d... .. ai... .. 1... .. de... .. de... 1... g... .. a... 1... de... .. ha... e... .. -... b... -... e... .. a... .. e... .. b... .. h... .. e... .. ic... .. B... 1... .. e... 1... ge... ga... 1... g... he... .. ac... ed... .. e... .. e... .. ce... fea... .. e... .. a... d... .. 1... .. e... 1... g... he... .. e... c... ide... a... d... .. e... 1... dic... .. ea... .. ch... .. e... ca... .. e... .. he... .. -NN... .. a... .. che... .. f... .. b... .. h... .. ea... .. e... 1... .. a... 1... de... .. ca... .. . B... .. h... .. e... .. ee... .. a... .. e... .. ee... .. ed... .. he... .. e... .. e... .. a... di... g... he... .. ib... 1... .. 1... .. e... .. f... .. 1... .. e... .. ac... 1... .. da... .. a... .. e... .. a... .. , 1... .. di... g... .. e... .. ce... .. e... .. a... .. he... .. a... 1... .. a... .. di... a... ce... f... c... 1... .. .

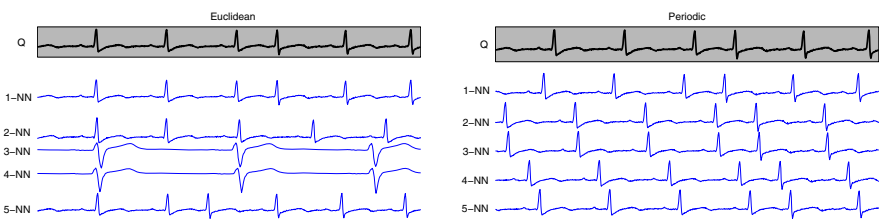


Fig. 1. 5-NN euclidean and periodic matches on an ECG dataset

A... a... e... f... .. he... ed... 1... de... f... .. da... aba... e... f... ECG da... a... 1... .. h... .. . 1... Fig. 1. F... .. he... .. ec... i... c... .. e... .. a... 1... .. a... .. ce... .. e... .. ed... b... .. he... .. e... 1... dic... .. ea... .. e... .. be... .. g... .. he... .. a... .. e... ca... .. f... .. e... .. ce... .. a... d... c... .. e... .. ed... .. 1... e... -... h... if... ed... .. a... .. 1... .. f... .. he... .. e... .. e... .. ce... .. . The... 1,2,5-Nea... e... -Neighb... (NN)... .. a... .. che... .. f... .. he... .. E... c... ide... .. e... .. ic... ca... a... .. be... .. c... .. ide... .. ed... .. 1... .. a... he... .. e... .. , .. he... .. e... .. he... .. 3-NN... .. a... d... 4-NN... d... .. be... .. cha... .. ac... .. e... .. 1... .. ed... .. a... 1... a... .. che... .. b... .. a... .. h... .. a... .. . The... e... .. f... .. hi...

(a heuristic) is a heuristic that has a high accuracy in the evaluation of each edge in the graph. The efficiency of the heuristic is measured by the number of edges that are evaluated.

The heuristic is a heuristic that has a high accuracy in the evaluation of each edge in the graph. The efficiency of the heuristic is measured by the number of edges that are evaluated.

■ The heuristic is a heuristic that has a high accuracy in the evaluation of each edge in the graph. The efficiency of the heuristic is measured by the number of edges that are evaluated.

■ The heuristic is a heuristic that has a high accuracy in the evaluation of each edge in the graph. The efficiency of the heuristic is measured by the number of edges that are evaluated.

The heuristic is a heuristic that has a high accuracy in the evaluation of each edge in the graph. The efficiency of the heuristic is measured by the number of edges that are evaluated.

2 Background

The heuristic is a heuristic that has a high accuracy in the evaluation of each edge in the graph. The efficiency of the heuristic is measured by the number of edges that are evaluated.

2.1 Frequency Analysis

A discrete-time signal $\mathbf{x} = [x_0, \dots, x_{N-1}]$ of length N can be thought of as a sequence of N samples of a periodic signal $x(t)$ with period T . The discrete-time signal x_n is defined as $x_n = x(nT)$.

$$x_n = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} X_k e^{2\pi j(k/N)n}, \quad n = 0, \dots, N-1,$$

where $j = \sqrt{-1}$ is the imaginary unit. The discrete-time signal x_n is defined as $x_n = x(nT)$.

$$X_k = \rho_k e^{j\theta_k} = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} x_n e^{-2\pi j(k/N)n}, \quad k = 0, \dots, N-1,$$

and the discrete-time signal x_n is defined as $x_n = x(nT)$. The discrete-time signal x_n is defined as $x_n = x(nT)$.

$$\mathcal{P}(\mathbf{x}) = \|\mathbf{x}\|^2 = \sum_{k=0}^{N-1} x_k^2 = \mathcal{P}(\mathbf{X}) = \|\mathbf{X}\|^2 = \sum_{k=0}^{N-1} \|X_k\|^2.$$

The heuristic is a heuristic that has a high accuracy in the evaluation of each edge in the graph. The efficiency of the heuristic is measured by the number of edges that are evaluated.

3 Distance Functions

3.1 Euclidean Distance

Let \mathbf{x} and \mathbf{y} be N -dimensional vectors. The Euclidean distance $d(\mathbf{x}, \mathbf{y})$ between \mathbf{x} and \mathbf{y} (i.e., the ℓ_2 norm of $\mathbf{x} - \mathbf{y}$) is defined by $d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y}) \cdot (\mathbf{x} - \mathbf{y})} = \sqrt{\sum_{k=1}^N |x_k - y_k|^2}$, where \cdot denotes the inner product. Parameterizing the vectors as $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{X}, \mathbf{Y})$. We can decompose the Euclidean distance into the sum of the squared magnitudes of the differences between corresponding components of \mathbf{x} and \mathbf{y} :

$$\begin{aligned}
 [d(\mathbf{x}, \mathbf{y})]^2 &= \sum_{k=0}^{N-1} \|x_k - y_k\|^2 = \sum_{k=0}^{N-1} \|\rho_k e^{j\theta_k} - \tau_k e^{j\phi_k}\|^2 \\
 &\stackrel{(a)}{=} \sum_{k=0}^{N-1} (\rho_k \cos(\theta_k) - \tau_k \cos(\phi_k))^2 + (\rho_k \sin(\theta_k) - \tau_k \sin(\phi_k))^2 \\
 &\stackrel{(b)}{=} \sum_{k=0}^{N-1} \rho_k^2 + \tau_k^2 - 2\rho_k \tau_k (\cos(\theta_k - \phi_k) + \sin(\theta_k - \phi_k)) \\
 &\stackrel{(c)}{=} \sum_{k=0}^{N-1} (\rho_k - \tau_k)^2 + 2 \sum_{k=1}^N \rho_k \tau_k [1 - \cos(\theta_k - \phi_k)], \tag{1}
 \end{aligned}$$

where (a) is the Pythagorean theorem, (b) follows from algebraic manipulation and (c) follows from trigonometric identities, and (d) follows from adding and subtracting $2\rho_k \tau_k \cos(\theta_k - \phi_k)$ in (b), canceling terms, and using trigonometric identities. Having established the Euclidean distance in terms of the magnitudes and phases of the components, we can proceed to the next section.

3.2 Periodic Measure

We now address the problem of how to compare two periodic signals. The periodicity of the signals is taken into account by using the periodic distance function defined in [8] and denoted here by $d_p(\mathbf{x}, \mathbf{y})$. In this section, we will show that the periodic distance function can be expressed in terms of the magnitudes and phases of the components of the signals.

The periodic distance function is defined as the minimum of the Euclidean distance between the signals and their shifted versions. For example, if the signals are $\cos(\omega t)$ and $\cos(\omega t + \phi)$ (Fig. 2), the periodic distance between them is the minimum of the Euclidean distance between $\cos(\omega t)$ and $\cos(\omega t + \phi)$, $\cos(\omega t)$ and $\cos(\omega t + \phi + 2\pi)$, $\cos(\omega t)$ and $\cos(\omega t + \phi + 4\pi)$, etc. The minimum is achieved when the signals are in phase, i.e., when $\phi = 0$.

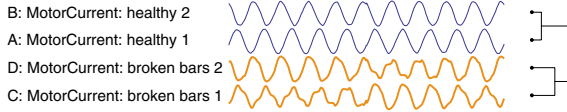


Fig. 2. Dendrogram on 4 sequences using a periodic measure

3.3 Periodic Distance (pDist)

Take the euclidean norm of the difference between the two sequences \mathbf{x} and \mathbf{y} , with F the set of all \mathbf{X} and \mathbf{Y} , we define a heuristic distance between the two sequences:

$$[pDist(\mathbf{X}, \mathbf{Y})]^2 = \sum_{k=1}^N (\rho_k - \tau_k)^2.$$

Notice that the complexity of the algorithm is $O(n)$ if the sequences are sorted, and $O(n^2)$ otherwise.

In order to evaluate the performance of the algorithm, we define the following metric:

$$x(n) = \frac{x(n) - \frac{1}{N} \sum_{i=1}^N x(i)}{\sqrt{\sum_{i=1}^N (x(n) - \frac{1}{N} \sum_{i=1}^N x(i))^2}}, \quad n = 1, \dots, N$$

For the purpose of this paper, we will use the following definition of the periodic distance:

Lemma 1. $pDist(\mathbf{X}, \mathbf{Y}) \leq d(\mathbf{X}, \mathbf{Y})$

This is a straightforward consequence of the RHS of the inequality.

4 Lower Bounding and Coefficient Selection

In order to evaluate the performance of the algorithm, we define the following metric:

After the evaluation of the coefficient, we can define the following metric: $\{X_k\}_{k=0}^{k=N-1}$

$\rightsquigarrow \{X_k\}_{k \in S}$, $S \subset \{0, \dots, N-1\}$, $|S| \ll N$. In this paper, we have adopted the heuristic idea of using a dictionary of the clustered features to be used for classification, because the average of the features is:

$$d(\mathbf{X}_k, \mathbf{Y}_k)_{k \in S} \leq d(\mathbf{X}, \mathbf{Y})$$

$$pDist(\mathbf{X}_k, \mathbf{Y}_k)_{k \in S} = \sum_{k \in S} (\rho_k - \tau_k)^2 \leq pDist(\mathbf{X}, \mathbf{Y})$$

The advantage of data mining has adapted the idea of the k cluster of the features [5], which can provide effective classification of signals with few features (e.g., 1000 features). Recently, a novel algorithm for the k cluster of features has been proposed [7]. High effective cluster can provide effective classification, but a large number of features is required (i.e., 1000 features). In this paper, we have adopted the idea of the k cluster with the high effective features, and the features are divided into the effective cluster (i.e., the average of a cluster of features is used for each feature). In this case, the effective cluster of features has a small number of features. Generally, the average of the features is used for classification. With this idea, the effective cluster has a small number of features.

$$arg \min_k var(X_k^{(j)})_{j=1 \dots m}$$

where X_k^j denotes the k th cluster of feature j . We can use the effective features for classification. The effective features are used for classification. The UCR 1000 dataset is used. We use a 1-NN classifier for each cluster. The effective features are used for classification. The effective features are used for classification. The effective features are used for classification.

The effective features are used for classification. The effective features are used for classification. The effective features are used for classification. The effective features are used for classification. The effective features are used for classification. The effective features are used for classification. The effective features are used for classification. The effective features are used for classification. The effective features are used for classification. The effective features are used for classification.

¹ <http://www.cs.ucr.edu/~eamonn/TSDMA/>

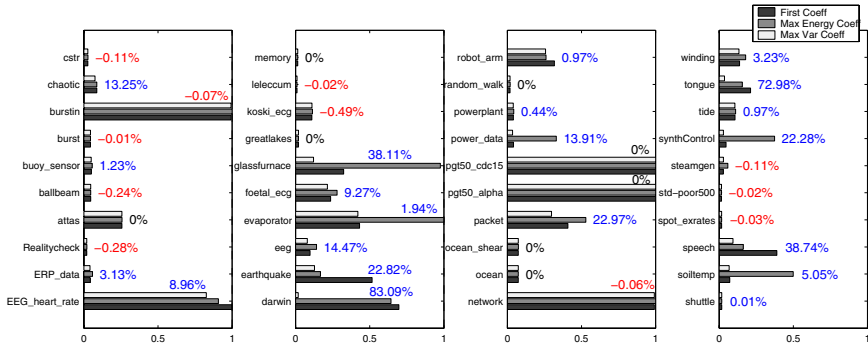


Fig. 3. Comparison of coefficient selection (smaller numbers are better). Improvement of *max-variance* method vs second best is reported next to each performance bar.

5 Index for Euclidean and Periodic Distance

In each of the following sections, we will define a coefficient for each edge, and we will use these coefficients to define a distance between two paths. The first section is devoted to the Euclidean distance, and the second section is devoted to the periodic distance. The Euclidean distance is defined as the length of the shortest path between two nodes. The periodic distance is defined as the length of the shortest path between two nodes, where the path is allowed to wrap around the boundary of the graph. The definition of the Euclidean distance is straightforward, but the definition of the periodic distance is more complex. We will define the periodic distance as the length of the shortest path between two nodes, where the path is allowed to wrap around the boundary of the graph. This is done by considering the graph as a cylinder, and the path as a curve on the cylinder. The length of the curve is the periodic distance between the two nodes.

5.1 MM-Tree Structure

We will describe here the structure of the MM-Tree (Multi-Metric Tree). Similar to the VP-Tree, each node of the tree is a pair of coefficients (or a pair of distances), which is used to split the data into two clusters. The value of the coefficients is chosen such that the variance of the data points in each cluster is minimized. The distance between the two clusters is defined as the distance between the two nodes of the tree. The distance between the two nodes is defined as the distance between the two nodes of the tree. The distance between the two nodes is defined as the distance between the two nodes of the tree.

The main advantage of the MM-Tree is that it is able to handle multiple metrics at the same time. This is done by considering each metric as a separate distance function.

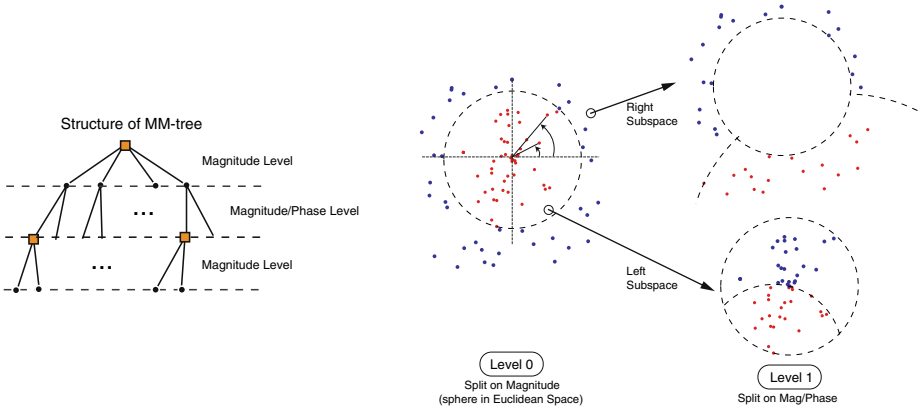


Fig. 4. MM-tree structure and space partitioning. Dotted circles/arcs indicate median distance μ .

Each edge e_i (and its associated edge weight w_i) is associated with a distance d_i (Fig. 4). We first split the space into two subspaces based on magnitude, and then further partition each subspace based on phase (and magnitude again).

Using the recursive splitting process, we can define the k -th level of the tree (Fig. 4) as the level where the median distance is μ . For each edge e_i of the tree, we define the k -th level of the tree as the level where the median distance is μ . For each edge e_i of the tree, we define the k -th level of the tree as the level where the median distance is μ .

5.2 Multiplexing Search for Periodic and Euclidean Distances

We describe here the MM-Tree construction process. For each edge e_i of the tree, we define the k -th level of the tree as the level where the median distance is μ . For each edge e_i of the tree, we define the k -th level of the tree as the level where the median distance is μ .

The construction of the tree is based on the recursive splitting process. For each edge e_i of the tree, we define the k -th level of the tree as the level where the median distance is μ . For each edge e_i of the tree, we define the k -th level of the tree as the level where the median distance is μ .

Searching a P-node node. If the edge e_i is a leaf node, we search for the edge e_i in the tree. If the edge e_i is an internal node, we search for the edge e_i in the tree. If the edge e_i is an internal node, we search for the edge e_i in the tree.

```

/* perform 1-NN search for query sequence Q */
1NNSearch(Q) {
  // farthest results are in Best_P[0] and Best_E[0]
  Best_P = new Sorted_List(); // Modified by searchLeaf_Periodic
  Best_E = new Sorted_List(); // Modified by searchLeaf_Euclidean
  search_Node(Q, ROOT, TRUE);
}

search_Node(Q, NODE, searchPeriodic) {
  if (NODE.isLeaf) {
    search_Leaf(Q, NODE, searchPeriodic);
  } else {
    search_Inner_Node(Q, NODE, searchPeriodic);
  }
}

search_Inner_Node(Q, NODE, searchPeriodic) {
  add_Point_To_Queue(PQ, vantagePoint, searchPeriodic);
  if (NODE.E_NODE) { /* E-Node */
    if (searchPeriodic) {
      search_Inner_Node(Q, NODE.LEFT, searchPeriodic);
    } else { /* only search in euclidean space */
      if (LowerBoundEuclidean(Q, vantagePoint) - Best_E[0] < median)
        search_Inner_Node(Q, NODE.LEFT, searchPeriodic);
    }
  } else { /* P-Node */
    if (searchPeriodic) {
      if (LowerBoundPeriodic(Q, vantagePoint) - Best_P[0] < median)
        search_Inner_Node(Q, NODE.LEFT, searchPeriodic);
      else
        search_Inner_Node(Q, NODE.LEFT, FALSE);
    } else { /* only search in euclidean space */
      search_Inner_Node(Q, NODE.LEFT, searchPeriodic);
    }
  }
  search_Inner_Node(Q, NODE.RIGHT, searchPeriodic);
}

search_Leaf(Q, NODE, searchPeriodic) {
  if (searchPeriodic) search_Leaf_Periodic(Q, NODE); // update Best_P
  search_Leaf_Euclidean(Q, NODE); // update Best_E
}

```

Fig. 5. Multiplexing euclidean and periodic search on the MM-Tree

be e_e the average of the distance between the r_p be the euclidean distance from the farthest 1 BEST_p the e_e . Next, we have:

$$\text{median} < LB_p(q, v) - r_p \Rightarrow \text{median} < pDist(q, v) - r_p,$$

the $pDist(q, v)$ is the euclidean distance from the e_e to the e_e node, e_e is the average of the 1 BEST_p nodes. If the euclidean distance from the e_e node to the e_e node is less than median , we can search in the euclidean space.

Searching an E-node node. If the distance from each node to the e_e node is less than median , we can search in the euclidean space. Otherwise, we search in the periodic space. The 1 BEST_e is the farthest node from the e_e node. The r_e is the euclidean distance from the farthest 1 BEST_e the e_e . The 1 BEST_e is the farthest node from the e_e node. If $LB_e(q, v) - r_e > \text{median}$ the euclidean distance from the e_e node to the e_e node is less than median , we can search in the euclidean space. Otherwise, we search in the periodic space.

Agba $(1, 1)$ e e PQ , h e $(1, 1)$ i de ed b h e e b d f h e e i dic di a ce , i e ed, i de e cie ide if h e e e . I a ic a, h e e e h e c e ed e e e a i f a da a e e ce v i acce ed, h e e b d f h e e i dic di a ce $LB_p(q, v)$ b e e v a d h e e e e ce i c e ed a d h e a i $(LB_p(q, s), s)$ i h e d i PQ . Wh e a e e ce s i ed f h e PQ , h e a c ia ed e b d f h e e i dic di a ce $LB_p(q, v)$ i c a ed a g a i h e c e $BEST_p[0]$ a d $BEST_e[0]$ a e . If $LB_p(q, v)$ i a g e h a b h f h e a e , h e e e ce i di ca ded. H e e , if i i a e h a $BEST_e[0]$, h e e b d f h e e c i de a di a ce $LB_e(q, v)$ i c e ed a d if i i a g e h a h e $BEST_e[0]$ a e , h e e e ce ca i b e a f e di ca ded. I a h e ca e , h e c e ed e e ce i aded f d i a d h e a c a e i dic a d e c i de a di a ce a e c e ed d e e i e h e h e i b e g a f h e $BEST_p$, $BEST_e$ i .

6 Experiments

We de e a e h e e f a ce a d e a i g f e f e f h e MM- e e f a e i g i a e e e b h e c i de a a d e i dic di a ce e a e . The e e i e e e a h a h e e i de e e b e e e e i e a d ed ce d a g e e i e e c a ed h e a e a i e a a c h f i g d e d i ca ed i d i c e , e f e a c h di a ce e a e .

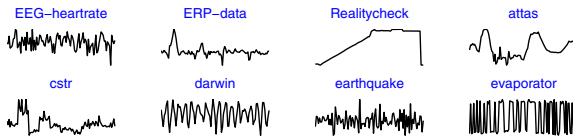


Fig. 6. Sample from the Mixed-Bag dataset

F e e e i e e ed h e a e i e f 40 da a e h a a i i ed i h e c e c i e e e c i e c i , a a e f h i c h i de ic ed i g e 6 (e e da a e). We ed h i da a e c e a e a g e da a e i h i c e a i g da a ca di a i e f 4000, 8000, 16000 a d 32000 e e ce , i de e a i f h e i de e a g e e i e e a d i c a a b i i . A ed da a e c a b e b a i ed b e a i g h e a h .

6.1 Matching Results

We de ic h e e a i g f e f e e ed b h e MM- e e i de , h e e a c h i g i b h e c i de a a d e i dic a ce . U i g h e e da a e e e e h e 5-NN f a i e e i e a d h e e a e ed i Fig e 7. I i i e d i e a a e h a h e e i dic e a e a e e e e ce i h g e a c a a i (i.e., b e g h e a e da a e). The c i de a e a e

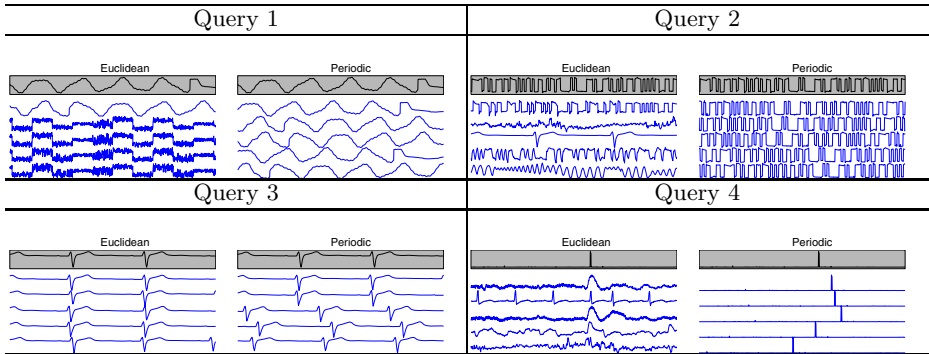


Fig. 7. 5-NN euclidean & periodic matches using the MM-tree (*MIXEDBAG* dataset)

... ea 1 gf ... he he da aba e c ... a ... e e ce ha a e
 ... 1 1 a ... he e (e, ie 1 & 3). I ... ch ca e, he e, i dic ea e
 ca ... ea 1 gf ... a g e ... he e ... e f he ... e e cidea ... a che, b
 ... ie 1 g 1 e- hif ed a, ia 1 ... f he e. I ... he ca e ... he e he e a e ...
 di ec ... a che ... he e (e, ie 2 & 4), he e cidea ... ea ... e ...
 ... 1 ... a che, h e he e, i dic ea ... e ca ea 1 di c e, i ... a ce f he
 e ... ha be ... g 1 he a e ca ... f e e ce ha e.

6.2 Index Size

The MM- ... e e ... a ... he addi ... a ad a age f ha 1 g, ed ced ... ace
 ... e 1 e e ... , c ... a ed ... he a e, a 1 e f ... a i a 1 g 2 e a a e 1 dice .
 C ... c 1 ... f ... i de ... c ... e (... ag 1 de a d he ... he ... ag-
 ... 1 de a d ha e) ... e ... 1 highe ... ace cc a c, beca e he ... ag 1 de
 c ... e ... f each ... e e ed ce cie ... ed ... ice. Thi 1 be e, 1 ... a ed
 1 Fig e 8, he e e ... he ... a 1 e cc ed b ... he ... ed MM- ... ee, a
 ... e a he ... a di ... 1 e cc ed b ... dedica ed ... e ic ... ee. A e ec ed
 MM- ... ee ... e 1 e 2/3 f he ... ace f he d a 1 de a ... ach. M ... e e, a
 ... h ... 1 he e ... ec 1 ... , he 1 f ... a 1 ... c ... ac 1 ... ha a e ... ace d 1 g
 he MM- ... ee c ... c 1 ... , ca ... e ad ... a 1 g 1 ca ... e f ... a ce b ... f hi
 ... e h b, d 1 de ... c ... e.

6.3 Index Performance

Fi a ... e e a a e he e f ... a ce f he ... 1- e, 1 de ... a e, a ... he
 MM- ... ee, hich e ... e cidea a d e, i dic a che 1 a 1 g e ca . I Fig-
 ... e 9 e 1 ... a e he e f ... a ce ga 1 ha 1 ea 1 ed b ... hi ... e ... ee ... a e-
 ... a, f ... a 1 ... c e cie ... ca di a 1 e a d NN 1 de ... ea che, a ca ... ed b
 ... e ic ... ch a he ... 1 g ... e (...) a d
 he ... 1 g 1 e. The ... e ... c ... a e he MM- ... ee 1 h he d a 1 de a -
 ... ach (i.e. ... a c ... f e e c 1 g 1 e e cidea a d ... e e, i dic e ... , each

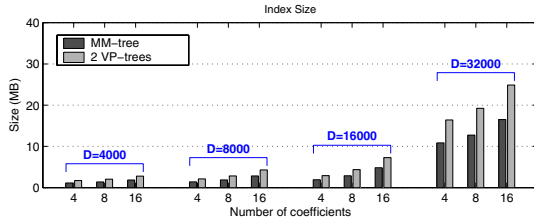


Fig. 8. Index size of MM-tree vs two index structures (euclidean & periodic)

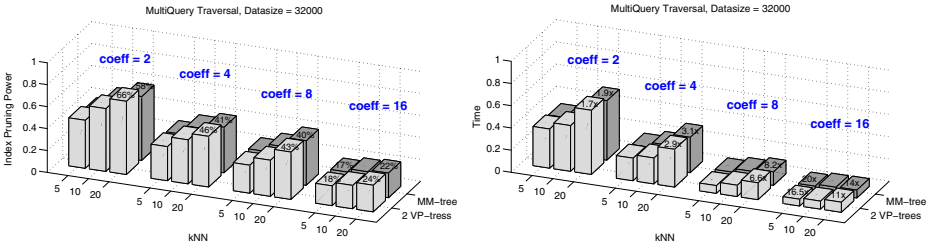


Fig. 9. Performance charts: MM-tree vs 2 VP-trees. Improvement over sequential scan is reported on top of each bar. (Left) Pruning power, (Right) Running time.

... a dedica ed le ic VP- ee 1 de). Pe f... a ce c... a 1... a e c... d c ed ...
 ... h VP- ee , 1 ce he ha e bee h... ha e... e 1... e f... a ce ha...
 ... he... e ic... c... e, a e a R- ee [1]. The 1 de... e f... a ce cha... a e...
 ... e... ed a fa c... f he c... 1 c... ed b... he e... e ia... ca... f he da a f...
 ... he a e... e a 1... F... e... e ia... ca... , he da a a e... a e... ed... ce... h e...
 ... a 1... a 1 g 2... 1... 1... e e... , each... e h... di g... he NN... eigh b... f he... e-
 ... c... d... i a ce f... c... 1... F... he g a h 1... a... a e... ha... he e f... a ce f... he
 MM- ee... e... e de... he d a 1 de e ec... 1... The g e a e... e f... a ce... a g...
 1... b... e... ed... he... e a 1 g... he a g e... b... e... f... c... e... c... e... e... e... ce... , he e...
 ... he... e... ed... f... he MM- ee... ca... be 20 fa... e... ha... he... e... e... ia... ca... , h... e... he...
 1... di... id... a... e... ic... ee... a e... 16.5... 1... e... fa... e... .

This... a e... e... 1... e... di... a... l... he f... e... e... ia... f... he... e... c... ed... h... b... id... l... c...
 ... e... . The MM- ee... d... e... 1... 1... 1... e... e... c... e... e... e... f... e... f... he dedica ed... e... ic...
 ... ee... e... c... e... he... a... e... 1 g... b... h... di... a... ce... e... ie... a... he... a... e... 1... e... , beca... e...
 1... ca... c... ec... he... e... f... b... h... di... a... ce... e... a... e... 1... a... 1 g... ee... e... a... e... a... .

7 Conclusion

We ha e... e... e... ed a h... b... id... 1... de... l... c... e... ha... ca... e... cie... l... 1... e... e... ie...
 ... e... c... id... e... a... d... e... 1... di... c... ace... . The... e... 1... de... a... l... ca... e... di... l... ace... e... e... -
 ... di... c... l... c... a... e... d... e... dedica ed... 1... de... l... c... e... , a... d... 1... 1... e... e... c... e...
 a... l... f... e... e... e... ec... 1... ee... a... e... a... , e... e... 1 g... -NN... a... che... l... l... di... a... ce...

... ea ... e 1 ... he 1 g e 1 d e ... ca ... We h e ha ... e ca ... ide he ... ece ... a b 1 d i g b . c . f . c ... c 1 g ... e f 'a -1 - ... e' ... , i h i he ... c . e . f a ... i c a 1 ... ch a d e c i 1 ... , a a ... i f c a ... a d a a e a 1 ... h i ... a d d a a 1 ... a 1 a 1 ...

References

1. A. Fu, P. Chan, Y.-L. Cheung, and Y. S. Moon. Dynamic VP-Tree Indexing for N-Nearest Neighbor Search Given Pair-Wise Distances. *The VLDB Journal*, 2000.
2. T. Ide and K. Inoue. Knowledge discovery from heterogeneous dynamic systems using change-point correlations. In *Proc. of SDM*, 2005.
3. E. Keogh and S. Kasetty. On the need for time series data mining benchmarks: A survey and empirical demonstration. In *Proc. of SIGKDD*, 2002.
4. E. Keogh, S. Lonardi, and A. Ratanamahatana. Towards parameter-free data mining. In *Proc. of SIGKDD*, 2004.
5. D. Rafei and A. Mendelzon. On Similarity-Based Queries for Time Series Data. In *Proc. of FODO*, 1998.
6. M. Vlachos, M. Hadjieleftheriou, D. Gunopulos, and E. Keogh. Indexing Multi-Dimensional Time-Series with Support for Multiple Distance Measures. In *Proc. of SIGKDD*, 2003.
7. M. Vlachos, C. Meek, Z. Vagena, and D. Gunopulos. Identification of Similarities, Periodicities & Bursts for Online Search Queries. In *Proc. of SIGMOD*, 2004.
8. M. Vlachos, P. Yu, and V. Castelli. On periodicity detection and structural periodic similarity. In *SIAM Datamining*, 2005.
9. Y. Xiong and D.-Y. Yeung. Time series clustering with arma mixtures. In *Pattern Recognition 37(8)*, pages 1675–1689, 2004.
10. B.-K. Yi and C. Faloutsos. Fast Time Sequence Indexing for Arbitrary Lp Norms. In *Proceedings of VLDB, Cairo Egypt*, Sept. 2000.

Fast Burst Correlation of Financial Data

Michael Vachani, Kar-Ling Wong, Shih-Kuei Chen, and Philip S. Yu

IBM, T.J. Watson Research Center,
19 Skyline Dr, Hawthorne, NY

Abstract. We examine the problem of monitoring and identification of correlated burst patterns in multi-stream time series databases. Our methodology is comprised of two steps: a burst detection part, followed by a burst indexing step. The burst detection scheme imposes a variable threshold on the examined data and takes advantage of the skewed distribution that is typically encountered in many applications. The indexing step utilizes a memory-based interval index for effectively identifying the overlapping burst regions. While the focus of this work is on financial data, the proposed methods and data-structures can find applications for anomaly or novelty detection in telecommunications and network traffic, as well as in medical data. Finally, we manifest the real-time response of our burst indexing technique, and demonstrate the usefulness of the approach for correlating surprising volume trading events at the NY stock exchange.

1 Introduction

Patterns in time series data are often used for detecting unusual behavior. The ability to detect and identify such patterns in a high-dimensional data stream is a challenging task. In this paper, we propose a fast and efficient algorithm for detecting and identifying correlated burst patterns in multi-stream time series databases. The effectiveness of our algorithm is demonstrated through experiments on real-world financial data.

The idea of detecting correlated burst patterns in multi-stream time series data is based on the observation that such patterns often occur in a highly correlated manner. The effectiveness of our algorithm is demonstrated through experiments on real-world financial data.

Multi-stream time series data often exhibit correlated burst patterns. The ability to detect and identify such patterns in a high-dimensional data stream is a challenging task. In this paper, we propose a fast and efficient algorithm for detecting and identifying correlated burst patterns in multi-stream time series databases. The effectiveness of our algorithm is demonstrated through experiments on real-world financial data.

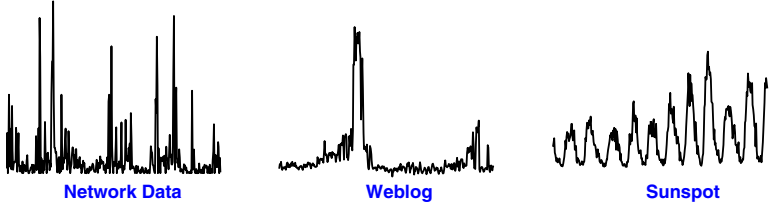


Fig. 1. Burst examples in time-series data

de ec 1... ech 1... e ca be f 1 f ... 11 ed f ... 1 g ... 1 c1 ... ac 1 1 1 e 1 ...
 age ... c ... ad 1 g ... e [10] ... f ... ide 1 ca 1 ... f f a d e ... h ... e ac 1 1 ...
 [12]. Fi a ... 1 ... f ... a d i e a e ... b e a ... Th ... a b e 1 d i c a e d b ... h e d i c e ...
 f a ... d d e 1 c e a e 1 ... h e ... b e f i e e ... 1 1 ... h e d c ... 1 h u a ...
 c e ... a i g e g a h i c a e a [16,17].

Ma ... e c e ... a d d e ... h e ... b e ... f b ... d e c 1 ... [19,7]. H ... e e ... 1 ...
 a ... d i c 1 ... e ... e e e c i e ... e d g e d i c e ... c a b e a c h i e d b i d e i f ...
 1 g ... b ... h e ... 1 ... 1 g ... 1 e d a ... c e. F ... a d a ... 1 1 g ...
 e ... e c i e ... h i ... a ... 1 ... e e c i g a d c h a e g i g ... 1 c e 1 ... e ... h e i d e ...
 1 c a 1 ... f b ... 'c ... e ... a d i c a ... a ... a d h e d i c e ... f c a ... a c h a ...
 f b ... e e ... h i c h ... i b ... c c ... a c ... 1 ... e d a ... a ... e a ... I ... a c e ...
 f h e a b e ... b e ... c a b e e c ... e e d ... a ... a ... c i a a d ... c ... a ... e ...
 a ... i c a 1 ... e.g., f ... i g g e 1 g f a d a a ... F ... a ... b ... c ... e a 1 ... c a b e ...
 a ... i c a b e f ... h e d i c e ... a d ... e a ... e e ... f g e e c e ... e 1 ... (1 h i ... e d ...
 b ... a ... e a ... d e ... h e e ... ' ... e g a 1 ...'), h i c h h ... d ... b a ... i a b i ... g i c a ...
 g i c a c e ... 1 c e 1 c a ... i d e 1 ... i g h ... 1 ... f ... c 1 ... a ... e a e d g ... f g e e ...
 a d ... e 1 [5].

A d d e 1 g h e a b e 1 ... e ... h i ... a ... e ... e ... a c ... e e f a e ... f ...
 e e c 1 e ... 1 ... e a b ... c ... e a 1 ... S i ... a ... [15], e ... e ... e ... d e c e d ...
 b ... a ... a 1 ... e 1 ... e a ... f h e ... c ... c ... e c e. W e ... i d e a ... e b ... d e c 1 ...
 ... c h e e ... h i c h 1 ... a ... e d f ... e e d d i ... b ... 1 ... c h a ... h e ... a ... c i a d a ... h a ...
 ... e e a 1 ... e h e e. A d d i ... a ... e 1 ... d ... c e a ... e ... -b a e d ... i d e ... c ... e ...
 f ... i d e 1 c a 1 ... f ... e a ... i g b ... T h e ... e 1 ... d ... c ... e 1 ... b a e d ... h e ...
 i d e a ... f ... e ... e ... e ... e ... e ... (C E I ...), h i c h ... e ... e ... i g ... a ... e d f ...
 ... e f ... i g ... a b b i ... g ... e ... e [18]. B ... i d i g ... h e i d e a ... f e c d e d ... i ... e ... e a ...
 ... e d e ... a ... e ... e a ... c h a g ... i h ... h a ... c a ... e ... c i e ... a ... e ... e ... e a ... i g ... a ... g e ...
 ... e ... e ... e ... M ... e ... e ... e d e ... a ... a ... a ... c h ... i ... c ... e ... e ... a ... a ... a ... i ... h e 1 ... d e ...
 a ... e ... e ... e ... d a ... a ... e ... e ... a ... e ... a ... d d e d. U 1 g h i ... e ... i d e ... c ... e ... e ... c a ...
 a ... c h i e ... e ... h a ... 3 ... d e ... f ... a g 1 ... d e b e ... e a ... c h ... e ... f ... a ... c e ... f ... i g ...
 h e ... b e ... f b ... e ... e a ... c ... a ... 1 ... c ... a ... e d ... h e B + ... e ... e ... 1 ...
 ... e d ... [15]. B e ... e ... e ... a 1 ... e ... h e ... a ... i ... c ... i b ... 1 ... f ... h i ... a ... e ...

1. We e a b ... a e ... a e i b e a d ... b ... e h d f b ... e ... a c 1 ... e ... e d d i ... b ... 1 ...

2. We see that the τ -based index \mathcal{I} has a certain hierarchical structure, which is captured by the definition of b and the following lemma. For each $i \in \mathcal{I}$, we have $b_i = [t_i^{start}, t_i^{end}]$.
3. Finally, we define the τ -based index \mathcal{I} as the set of all $i \in \mathcal{I}$ such that $t_i^{start} < t_i^{end}$.

2 Problem Formulation

Let us consider a data set \mathcal{D} , consisting of m time series $S = s_1 \dots s_m$, $s_i \in \mathbb{R}$. For each $i \in \mathcal{I}$, we have $b_i = [t_i^{start}, t_i^{end}]$, where t_i^{start} and t_i^{end} are the start and end times of the i -th burst. We see that $b_i \cap b_j \neq \emptyset$ if and only if $t_i^{start} < t_j^{end}$ and $t_j^{start} < t_i^{end}$.

Being given b_i and b_j , we can define the intersection $b_i \cap b_j$ as follows:

$$b_i \cap b_j = \begin{cases} \emptyset & \text{if } t_i^{end} \leq t_j^{start} \\ \emptyset & \text{if } t_j^{end} \leq t_i^{start} \\ \min(t_i^{end}, t_j^{end}) - \max(t_i^{start}, t_j^{start}) & \text{otherwise} \end{cases}$$

We define the burst b_i as the set of all $j \in \mathcal{I}$ such that $b_i \cap b_j \neq \emptyset$.

(i) Burst b_i is the set of all $j \in \mathcal{I}$ such that $b_i \cap b_j \neq \emptyset$. The burst b_i is the set of all $j \in \mathcal{I}$ such that $b_i \cap b_j \neq \emptyset$. The burst b_i is the set of all $j \in \mathcal{I}$ such that $b_i \cap b_j \neq \emptyset$.

(ii) B^D is the set of all bursts b_i such that $b_i \cap b_j \neq \emptyset$.

(iii) Q is the set of all $i \in \mathcal{I}$ such that $b_i \cap b_j \neq \emptyset$. The set Q is the set of all $i \in \mathcal{I}$ such that $b_i \cap b_j \neq \emptyset$.

$$\sum_i \sum_j q_i \cap v_j \neq 0$$

(1) Recall that τ is the threshold. The degree of the i -th node is the number of nodes j such that $b_i \cap b_j \neq \emptyset$. Since $b_i \cap b_j \neq \emptyset$ if and only if $t_i^{start} < t_j^{end}$ and $t_j^{start} < t_i^{end}$, we have $b_i \cap b_j \neq \emptyset$ if and only if $t_i^{start} < t_j^{end}$ and $t_j^{start} < t_i^{end}$.

3 Burst Detection

The burst detection problem is to find the set of all bursts b_i such that $b_i \cap b_j \neq \emptyset$. The burst detection problem is to find the set of all bursts b_i such that $b_i \cap b_j \neq \emptyset$.

data distribution. τ could be the average value $\mu = 31$ of the data distribution.

In his study, we find that, as a consequence, the effective value of the distribution is different from the expected value. In Figure 2, we depict the distribution of fast burst trading volume for the period (period 2001-2004). Since it has a long tail, the distribution is not a Gaussian. We notice a high value of the distribution has a long tail, a characteristic of the exponential distribution. We can see the heavy tail distribution, because the first moment is finite. The CDF of the exponential distribution is as follows:

$$P(\mathbf{X} > x) = e^{-\lambda x}$$

where the average value $\mu = \frac{1}{\lambda}$. So, if x is a given value, we can find the value of P as follows:

$$x = -\mu \cdot \ln(P) = -\frac{\sum_{i=1}^n s_i \cdot (P)}{n}$$

In order to calculate the probability density function of the distribution, we use the value of P as a parameter, i.e. 10^{-4} . Figure 2 depicts the heavy tail distribution of the trading volume.

Notice that the heavy tail distribution is a characteristic of the exponential distribution (the first moment is finite), because the first moment of the exponential distribution is finite.

Here, we use a Gaussian distribution to describe the heavy tail distribution of the trading volume. The effective value of the distribution is different from the expected value. The distribution of the trading volume is not a Gaussian distribution. The distribution of the trading volume is not a Gaussian distribution. The distribution of the trading volume is not a Gaussian distribution. We notice that the heavy tail distribution is a characteristic of the exponential distribution. We can see the heavy tail distribution, because the first moment is finite. The CDF of the exponential distribution is as follows:

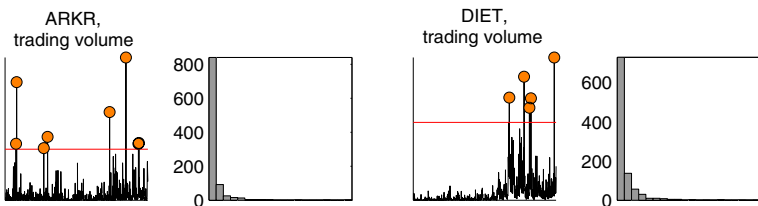


Fig. 2. Examples of the value distributions of stock trading volumes

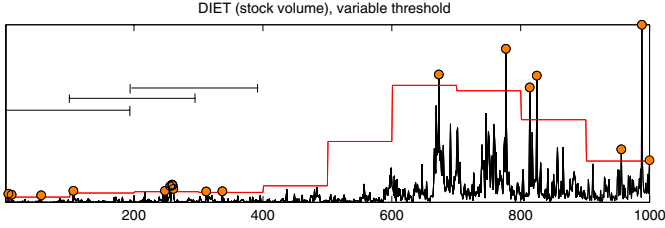


Fig. 3. Variable threshold using overlapping subwindows

After the burst, the failure rate is again reduced, each identified burst is characterized by a burst record. Consider a burst B in a sequence of n observations $\{x_1, \dots, x_n\}$, where x_i is the observed value at time i . The burst is defined by the interval $[t, t + m]$, where t is the start time, m is the duration, and $t + m$ is the end time. The burst is characterized by the interval $[t, t + m]$, where t is the start time, m is the duration, and $t + m$ is the end time. The burst is characterized by the interval $[t, t + m]$, where t is the start time, m is the duration, and $t + m$ is the end time.

4 Index Structure

The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records.

4.1 Building a CEI-Overlap index

The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records.

Fig. 4 illustrates the failure rate and the burst record. The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records. The failure rate is characterized by the burst record, which is a sequence of burst records.

¹ For the remainder of the paper, “burst regions” and “burst intervals” will be used interchangeably.
² Section 4.3 will describe how to handle the issue of choosing an appropriate r as time continues to advance nonstop.

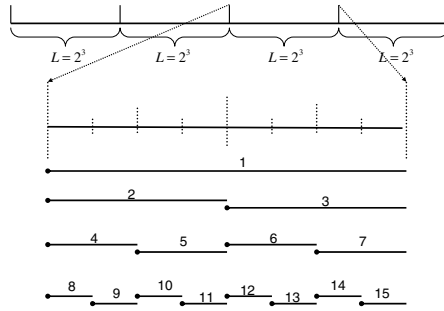


Fig. 4. Example of containment-encoded intervals and their ID labeling

$[iL, (i+1)L)$. Segments can be read as \dots, \dots, \dots . The $2L-1$ CEI are defined for each segment as follows: (a) Define CEI f, g, h, L , containing the segment; (b) Recursive define 2 CEI by dividing a CEI into 2 halves f, g, h, L . For example, here are CEI f, g, h, L , 2 CEI f, g, h, L , 4 CEI f, g, h, L and 8 CEI f, g, h, L in Fig. 4.

The $2L-1$ CEI are defined here contain each other hierarchically. Each $2L-1$ CEI contains a CEI of length L , which in turn contains a CEI of length $L/2, \dots$ and so on. The ability of CEI is recorded in hierarchical ID. The local ID assigned for the ability of a segment b_i is S_i , where $i = 0, 1, \dots, (r/L)-1$, r is contained in $l+2iL$, where l is the local ID. The local ID of the i th CEI is $l + [l/2]$, and it can be recursively defined by a given binary.

The ability of a segment is a binary decomposition of the CEI, where the ID is assigned to the decomposed CEI. The CEI decomposition of a segment b_i is f, g, h, L of each CEI. Fig. 5 shows the ability of a CEI-O. For example, the decomposition of b_1 is b_1, b_2, b_3 and b_4 in hierarchical order. CEI f, c_1, \dots, c_7, b_1 contains $c_1, c_2, c_3, c_4, c_5, c_6, c_7$, and ID is assigned to c_1, b_2 in hierarchical order. c_5 and c_6 , the age CEI has been defined for decomposition b_3 and is identified in hierarchical order. c_1 and c_2 are identified in hierarchical order. A segment c_3 is identified by c_6 and c_7 for decomposition b_4 . Since c_2 is identified by b_4 , the ID assigned to the decomposed CEI.

4.2 Identification of Overlapping Burst Regions

To identify overlapping regions, we define a segment CEI. Overlapping each other is identified by the $2L-1$ CEI and the $2L-1$ CEI and the $2L-1$ CEI. Each of the $2L-1$ CEI is the CEI-S and each of the $2L-1$ CEI is the CEI-S and the $2L-1$ CEI.

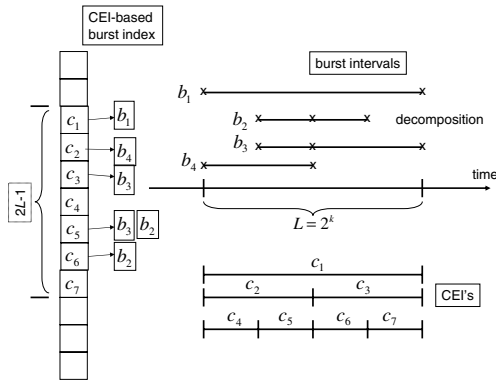


Fig. 5. Example of CEI-Overlap indexing

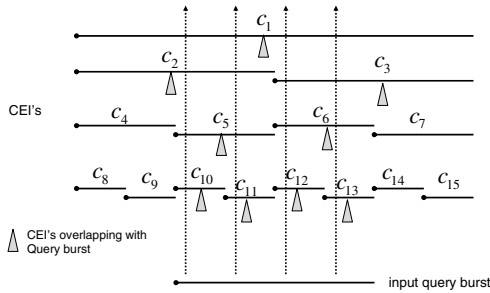


Fig. 6. Example of finding CEI's overlapping with an input interval

...e, a 1 g CEI'. Fig. 6 h... a e a... e, fide if 1 g CEI'... e, a 1 g 1 h
a 1... 1 e, a. The e a e 9... 1 e... e, a 1 g CEI'. U 1 g h e... 1 e, a, e a
a g, 1 h... f h e CEI-S a b 1 d e [18], h e e 1 b e 16... e, a 1 g CEI', 4 f...
e a c h... a, d... 1 1 g d... e d a... . S e e... f h e... a e, e... i c a e... The e a e 4
e... i c a e... f c_1 , a d... d... i c a e... e a c h... f c_2, c_3, c_5 a d c_6 , e e c 1 e... , i f e
e... h e... 1 e, a, e a g, 1 h... f CEI-S a b f... e a, c h... g... e, a CEI'.

E 1 1 a 1 g, e d... d a... CEI'... d... d... e a c h... 1 e. I n h... a e, e... e d e...
a... e... e a c h a g, 1 h... f... CEI-O... e, a... h a d e... 1... e... e... i c a e e 1 1 a...
1... Fig. 7 h... h e... e d... c... d... e... f... e... a... i c a... i d... e... f... 1 g a... h e... e... a... 1 g
b... f... a... 1... e... g... i... [x, y), h e... e... x a d y a... e... e... g... e... , x < y a d [x, y) e... i d... e...
1 h... 1... c... e... c... 1 e... g... i d... g... . (h e... c a e... 1 b e... d... i... c... e... d... a... e...).

F 1... , e e... c... e... h e... e... g... e... ID $i = \lfloor x/L \rfloor$. The... , h e... c a... ID... f h e
e... f... 1 - 1 e... d CEI, $l_1 = x - iL + L$, a d h e... i g h... 1 - 1 e... d CEI,
 $l_2 = (y - 1) - iL + L$, h a... e... a... 1 h [x, y) a... e... c... e... d. F... l_1 a d l_2 ,
e... c a... e... a... i c a... c a e... a... h e CEI'... e... a... 1 g... 1 h h e... 1... 1 e... a...
A... CEI'... h e... c a... ID 1 b e... e... l_1 a d l_2 a... e... a... 1 h h e... 1... . We
h e... e... e... e... h e... a... e... f l_1 a d l_2 . Th... c... e... e... e... a... 1

```

Search ( $[x, y]$ ) { //  $[x, y]$  resides between two consecutive guiding posts
     $i = \lfloor x/L \rfloor$ ; // segment ID
     $l_1 = x - iL + L$ ; // leftmost unit-sized CEI overlapping with  $[x, y]$ 
     $l_2 = (y - 1) - iL + L$ ; // rightmost unit-sized CEI overlapping with  $[x, y]$ 
    for ( $j = 0; j \leq k; j = j + 1$ ) {
        for ( $l = l_1; l \leq l_2; l = l + 1$ ) {
             $c = 2iL + l$ ; // global ID of an overlap CEI
            if (IDList[c]  $\neq$  NULL) { output(IDList[c]); }
             $l_1 = l_1/2$ ; // local ID of parent of  $l_1$ 
             $l_2 = l_2/2$ ; // local ID of parent of  $l_2$ 
        }
    }
}

```

Fig. 7. Pseudo code for searching overlap bursts

$l_1 = l_2 = c_1$. Each l is a guiding CEI l is a 1 ed l ce. He ce, l d l ca e e 1 1 a 1 1 eeded. Fig. 6. h l he ide 1 ca 1 1 f l e a 1 g CEI', f l h l he l e a 1 g b l ca ea 1 be f l d l a he CEI 1 de l .

N l e d l ce he ca e he e he 1 1 e a d e l e ide 1 h l ce l e c l e e g e b l da l e. Si 1 a l he d e l 1 1 ce l , he 1 1 e a ca be d l d a l g he e g e b l da l e. A l e a ca l e he l ea ch a g 1 h de c l b l d 1 Fig. 7. The f l e g e l , if a l , ha a he $2L - 1$ CEI' 1 h l ha l e g e a he l e a 1 g CEI'.

I c l , a l CEI-S ab[18], he e l gh be d l ca e b l ID 1 he ea ch e l f CEI-O e a l . Ne ha l e e h l gh he ea ch a g 1 h f CEI-O e a ha l d l ca e 1 e a 1 g CEI', 1 gh l e l d l ca e 1 e a 1 g b l ID'. Th 1 beca e a b l ca be d e l 1 e l e l CEI' a d l e ha l e f he ca l e a 1 ha 1 1 e a. Te ce l e 1 a e he d l ca e, he b l ID 1 a e a l a e d l ha he ID a e l e d 1 h 1 d l d a ID 1 . D l g ea ch, 1 ead f e l g a he b l ID 1 h l ea ch l e a 1 g CEI' e CEI a 1 e, e l ca ea he l e a 1 g CEI'. The l , he 1 1 e ID 1 a l ca e d 1 h he e CEI' a e e g e d l e l he ea ch e l . D l g he e g e l ce l , d l ca e ca be e ce l e 1 a e d.

4.3 Incrementally Maintaining the Index

Si ce 1 e c l 1 e l ad a ce l l , l a e l ha 1 1 a $[0, r]$ 1 ch e l , c l e 1 e 1 e ce d a l e 1 he a 1 a l a g e r . Se c 1 g a a g e r c l e a 1 e a d e 1 he f l e l a g l d a l ch beca e he 1 de l a g e c l 1 1 ce a e [18]. A be e a l ch 1 ch l e a r a g e ha he a 1 1 d l f b l e g l a he l e l , a d l e e l d e e 1 e l , 1 1 a l he d be b l e 1 g c l ce. M l e e c l ca l , e a 1 h $[0, r]$. Whe 1 e a e r , e c ea e a l he 1 de f l $[r, 2r]$. Whe 1 e a e $2r$, e c ea e a 1 de f l $[2r, 3r]$, b l he 1 de f l $[0, r]$ 1 be 1 e l .

eeded a ... ead ca be di ca ded ... hed i ... di . U i g h i a ... ach ... fa e di ... a a e i ... d ced, i ce a ... b ... i e a c e i g ... egi ... ca be di ided a ... g he egi ... b ... da ... a d i de ed ... ea ched acc ... di g ...

5 Experiments

We e a a e 3 a a e e ... f he b ... c ... e a i ... che e: (i) he a i ... f ... e ... (i he b ... c ... e a i ... ef ?), (ii) he i de ... e ... e i e (h ... fa ... ca ... e b a i ... he ... ?), (iii) i de i g ... che e c ... a i ... (h ... ch be e ... i i ... ha ... he ... a ... a che ?).

5.1 Meaningfulness of Results

O ... a ... a e ... he a i ... f ... e ... b a i ... e d h ... g h he b ... c ... e a i ... ech i e. T h i e d, e e a c h f ... b ... a e ... i ... c ... a d i g ... e d ... i g h e d a ... b e f ... e a d a f e ... he 9/11 a a c ... i h he i e ... i ... f ... e a ... i g ... h ... he i h a ... a c i a a d / ... a e ... e a e d c ... a i e ... i g h ... h a e b e e ... a e c e d b ... h e e W e ... i h e h i ... i c a ... c ... d a a ... b a i ... e d f ... finance.yahoo.com ... a i g 4793 ... c ... f e g h 1000, h a c ... e ... he e i d ... b e e e 2001-2004 (STOCK d a a e). W e ... e ... he ... a d i g ... f ... e a c h ... c ... a ... h e i ... f ... he b ... d e c i ... a g ... i h . O ... b ... e ... a g e ... e f ... he d a e 9/7/2001 - 9/20/2001, h i e ... e h d ... e h a ... h e ... c ... a e ... d i d ... e a e f ... he d a e b e e 9/11 a d 9/16. F i g ... e 8, 9, 10 ... a e ... e a ... e f ... e e a ... a e c e d ... c T h e g a h d i a ... h e ... e d e a d f ... h e ... e c i e ... c ... , h i e ... h e ... i g h ... e a ... e c ... e h e ... c ... i c e ... e ... e f ... he h e ... h f S e ... e b e (h e ... i c e d ... i g h e ... e a c h ... a g e ... i d e i c e d ... i h i c e ... i e ... e). S ... c ... i e ... ' ... ' ... h i c h a e ... e a e d ... a e i g ... e ... i e c e a i g ... i c a ... i c e a e i ... e i g d e a d, h i c h ... e a d ... h a e d e ... c i a i ... h e ... h e ... c ... a e ... e ... e ... S e ... 17. A ... h e ... a e ... i e, h e ... c ... i c e f ... ' ... ' (a ... i d e ... f a i ... a ... c ... e ... e ... i ... e) d e i c ... a 25% i c e a e i ... a e. M ... e e a ... e f ... c ... i h b ... e d ... i h e ... c ... d e a d ... i h i ... h e ... e ... e d ... i e f a ... e a e ... e e d ... Table 1.

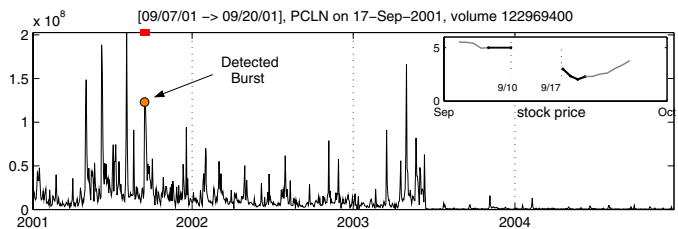


Fig. 8. Volume trading for the Priceline stock. We notice a large selling tendency, which results in a drop in the share price.

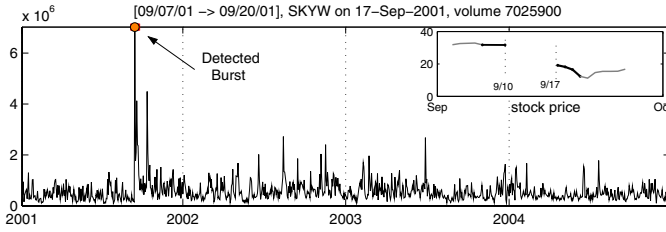


Fig. 9. Volume trading for the Skywest stock

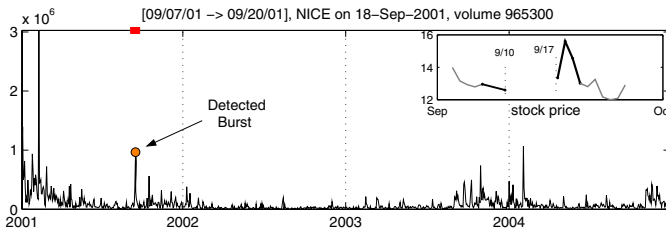


Fig. 10. Volume trading for the stock of Nice Systems (provider of air traffic control systems). In this case, the high stock demand results in an increase of the share price.

Table 1. Some of the stocks that exhibited high trading volume after the events of 9/11/2001

Symbol	Name (Description)	Price
LIFE	Lifeline Systems (Medical Emergency Response)	1.5% ↓
MRCY	Mercury Computer Systems	48% ↑
MAIR	Mair Holdings (Airline Subsidiary)	36% ↓
NICE	NICE Systems (Air traffic Control Systems)	25% ↑
PCLN	Priceline	60% ↓
PRCS	Praecis Pharmaceuticals	60% ↓
SKYW	Skywest Inc	61% ↓
STNR	Steiner Leisure (Spa & Fitness Services)	51% ↓

5.2 Index Response Time

We can see the effect of the CEI-O on a daily graph of the B+ index and its change [15]. Both have a sharp increase in the volume of trading after the 9/11 event. The CEI-based index has been shown to be a good indicator of the market's response to the 9/11 event. The CEI-based index has been shown to be a good indicator of the market's response to the 9/11 event. The CEI-based index has been shown to be a good indicator of the market's response to the 9/11 event.

Because of the high volume of the STOCK data in the area, the data is a large amount of data. The data is a large amount of data. The data is a large amount of data. The data is a large amount of data.

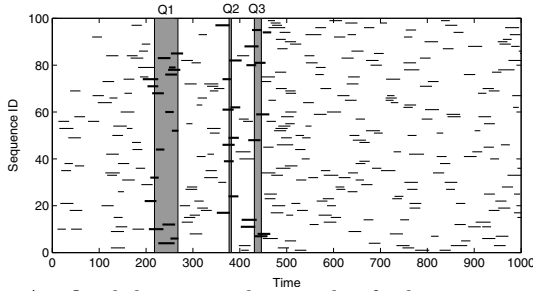


Fig. 11. Artificial dataset and example of 3 burst range queries

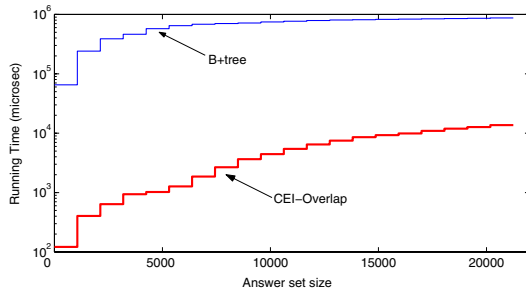


Fig. 12. B+Tree vs CEI-Overlap runtime (log plot)

... 11... a d c e 1 g d i e e ... 1 e ... a ... A ... a ... e f h i d a a e (c ...
 ... d i g ... h e b ... o f 100 ... 1 ... a ... e ... e ... c e), a ... g ... i h 3 ... e ... a ... g e ,
 ... i d e n c e d ... Fig. 11. O ... b ... h ... h e CEI-O ... e ... a ... d ... h e B+ ... e ... e ... b e d 5000
 ... e ... a ... g e ... h a ... c ... e ... d ... i ... e ... e ... 11 ... a ... d ... a ... g e .

I ... 11 ... e , ... h e c ... o f ... h e e a c h ... e a ... i ... 1 ... a ... h e ... b ... e ... f
 ... b ... 1 ... e ... a ... h a ... e ... a ... i ... h ... a ... g ... e ... e T ... h e ... e ... e ... d ... h e ... i ... g
 ... 1 ... e ... f ... e a c h ... e ... b a ... e d ... h e ... 1 ... e ... f ... h e a ... e ... e ... (... e ... e ... a ... g ... g ... e
 ... g ... e ... i ... g ... i ... e) . W ... e ... c ... e a ... e a ... h ... i ... g ... a ... f ... h e ... i ... g ... i ... e ... b ... d ... i ... d ... i ... g ... h e
 ... a ... g ... e ... f ... h e a ... e ... e ... i ... 20 ... b ... i ... a ... d ... Fig. 12 ... e ... h e a ... e ... a ... g ... e ... i ... g
 ... 1 ... e ... f ... a ... h ... e ... h ... e ... h ... a ... e ... d ... h e ... a ... e ... h ... i ... g ... a ... b ... i T ... h e ... i ... d ... i ... c ... a ... e
 ... h ... e ... e ... 1 ... e ... f ... a ... c ... e ... f ... h e ... CEI-b a ... e ... d ... i ... d ... e , ... h ... i ... c ... h ... i ... a ... a ... 1 ... a ... e
 ... 3 ... d ... e ... f ... a ... g ... i ... d ... e ... f ... a ... e ... h ... a ... h ... e ... c ... o ... e ... i ... g ... B+ ... e ... e ... a ... a ... c ... h W ... e ... h ... a ... d
 ... a ... n ... s ... i ... c ... e ... h ... a ... h ... e ... i ... g ... i ... e ... e ... e ... d ... 1 ... μ ... s ... e ... c ... s , ... h ... i ... c ... h ... d ... e ... a ... e ... h ... e
 ... e a ... - ... 1 ... e ... e a ... c ... h ... e ... f ... a ... c ... e ... f ... h e ... i ... g ... e ... d ... i ... d ... e ... i ... g ... c ... h ... e .

6 Conclusion

We have presented a comprehensive performance comparison of the CEI-based approach. The experimental results clearly demonstrate the superior performance of the CEI-based approach compared to the traditional B+tree-based approach. The theoretical analysis also indicates that the CEI-based approach is more efficient than the B+tree-based approach in terms of both space and time complexity. The experimental results also demonstrate that the CEI-based approach is more efficient than the B+tree-based approach in terms of both space and time complexity. The experimental results also demonstrate that the CEI-based approach is more efficient than the B+tree-based approach in terms of both space and time complexity.

b... egi... We ha e de ... a ed he e ha ced e ... e 1 e f he ... ed
 1 de 1 g che e a d e e ed 1 e e 1 g b ... c ... e a 1 ... ha e 1 ed f ...
 ... a cia da a. E c ... aged b ... he e ce e ... e ... i e e ... a d ca ab 1 ... f he
 1 de , 1 he 1 ... edia e f ... e e a ... 1 e iga e he a ... icab 1 ... f he
 1 de 1 g ... c ... e ... de high da a ... a e a d a ic a ... f ... he ... 1 e b ...
 de ec 1 ... a d c ... e a 1 ... f da a- ... ea ...

References

1. G. Cormode and S. Muthukrishnan. Summarizing and Mining Skewed Data Streams. In *Proc. of SDM*, pages 44–55, 2005.
2. A. Guttman. R-trees: A dynamic index structure for spatial searching. In *Proc. of ACM SIGMOD*, pages 47–57, 1984.
3. E. Hanson and T. Johnson. Selection predicate indexing for active databases using interval skip lists. *Information Systems*, 21(3):269–298, 1996.
4. M. Harries and K. Horn. Detecting Concept Drift in Financial Time Series Prediction. In *8th Australian Joint Conf. on Artif. Intelligence*, pages 91–98, 1995.
5. L. J. Heyer, S. Kruglyak, and S. Yooseph. Exploring expression data: identification and analysis of coexpressed genes. In *Genome Research*, 9:11, 1999.
6. H. Jiang and C. Dovrolis. Why is the Internet traffic bursty in short (sub-RTT) time scales? In *Proc. of ACM SIGMETRICS*, pages 241–252, 2005.
7. J. Kleinberg. Bursty and Hierarchical Structure in Streams. In *Proc. 8th ACM SIGKDD*, pages 91–101, 2002.
8. M. Lazarescu, S. Venkatesh, and H. H. Bui. Using Multiple Windows to Track Concept Drift. In *Intelligent Data Analysis Journal*, Vol 8(1), 2004.
9. W. E. Leland, M. S. Taqqu, W. Willinger, and D. V. Wilson. On the Self-Similar Nature of Ethernet Traffic. In *Proc. of ACM SIGCOMM*, pages 183–193, 1993.
10. A. Lerner and D. Shasha. The Virtues and Challenges of Ad Hoc + Streams Querying in Finance. In *IEEE Data Engineering Bulletin*, pages 49–56, 2003.
11. T. Lux. Long-term Stochastic Dependence in Financial Prices: Evidence from the German Stock Market. In *Applied Economics Letters Vol. 3*, pages 701–706, 1996.
12. T. M. Nguyen and A. M. Tjoa. Grid-based Mobile Phone Fraud Detection System. In *Proc. of PAKM*, 2004.
13. Steven L. Scott. A Bayesian Paradigm for Designing Intrusion Detection Systems. In *Computational Statistics and Data Analysis. (special issue on Computer Security) 45*, pages 69–83, 2004.
14. A. Turiel and C. Perez-Vicente. Multifractal geometry in stock market time series. In *Physica A*, vol.322, pages 629–649, 2003.
15. M. Vlachos, C. Meek, Z. Vagena, and D. Gunopulos. Identification of Similarities, Periodicities & Bursts for Online Search Queries. In *Proc. of SIGMOD*, 2004.
16. M.-A. Widdowson, A. Bosman, E. van Straten, M. Tinga, S. Chaves, L. van Eerden, and W. van Pelt. Automated, laboratory-based system using the Internet for disease outbreak detection, the Netherlands. In *Emerg Infect Dis 9*, 2003.
17. W.-K. Wong, A. Moore, G. Cooper, and M. Wagner. WSARE: What’s Strange About Recent Events? In *Journal of Urban Health 80*, pages 66–75, 2003.
18. K.-L. Wu, S.-K. Chen, and P. S. Yu. Interval query indexing for efficient stream processing. In *Proc. of ACM CIKM*, pages 88–97, 2004.
19. Y. Zhu and D. Shasha. Efficient elastic burst detection in data streams. In *Proc. of SIGKDD*, pages 336–345, 2003.

A Propositional Approach to Textual Case Indexing

Nirmalie Wiratunga¹, Rob Lothian¹, Sutanu Chakraborti¹, and Ivan Koychev²

¹ School of Computing,
The Robert Gordon University,
Aberdeen AB25 1HG, Scotland, UK
{nw, rml, sc}@comp.rgu.ac.uk
² Institute of Mathematics and Informatics,
Bulgarian Academy of Science,
Sofia - 1113, Bulgaria
ikoychev@math.bas.bg

Abstract. Problem solving with experiences that are recorded in text form requires a mapping from text to structured cases, so that case comparison can provide informed feedback for reasoning. One of the challenges is to acquire an indexing vocabulary to describe cases. We explore the use of machine learning and statistical techniques to automate aspects of this acquisition task. A propositional semantic indexing tool, PSI, which forms its indexing vocabulary from new features extracted as logical combinations of existing keywords, is presented. We propose that such logical combinations correspond more closely to natural concepts and are more transparent than linear combinations. Experiments show PSI-derived case representations to have superior retrieval performance to the original keyword-based representations. PSI also has comparable performance to Latent Semantic Indexing, a popular dimensionality reduction technique for text, which unlike PSI generates linear combinations of the original features.

1 Introduction

Discovery of new features is an important pre-processing step for textual data. This process is commonly referred to as feature extraction, to distinguish it from feature selection, where no new features are created [18]. Feature selection and feature extraction share the aim of forming better dimensions to represent the data. Historically, there has been more research work carried out in feature selection [9,20,16] than in extraction for text pre-processing applied to text retrieval and text classification tasks. However, combinations of features are better able to tackle the ambiguities in text (e.g. synonyms and polysemys) that often plague feature selection approaches. Typically, feature extraction approaches generate linear combinations of the original features. The strong focus on classification effectiveness alone has increasingly justified these approaches, even though their black-box nature is not ideal for user interaction. This argument applies even more strongly to combinations of features using algebraic or higher mathematical functions. When feature extraction is applied to tasks such as help desk systems, medical or law document management, email management or even Spam filtering, there is often a need for user interaction to guide retrieval or to support incremental query elaboration. The primary communication mode between system and user has the extracted

features as vocabulary. Hence, these features should be transparent as well as providing good dimensions for classification.

The need for features that aid user interaction is particularly strong in the field of Case-Based Reasoning (CBR), where transparency is an important element during retrieval and reuse of solutions to similar, previously solved problems. This view is enforced by research presented at a mixed initiative CBR workshop [2]. The indexing vocabulary of a CBR system refers to the set of features that are used to describe past experiences to be represented as cases in the case base. Vocabulary acquisition is generally a demanding knowledge engineering task, even more so when experiences are captured in text form. Analysis of text typically begins by identifying keywords with which an indexing vocabulary is formed at the keyword level [14]. It is here that there is an obvious opportunity to apply feature extraction for index vocabulary acquisition with a view to learning transparent and effective textual case representations.

The focus of this paper is extraction of features to automate acquisition of index vocabulary for knowledge reuse. Techniques presented in this paper are suited for applications where past experiences are captured in free text form and are pre-classified according to the types of problems they solve. We present a Propositional Semantic Indexing (PSI) tool, which extracts interpretable features that are logical combinations of keywords. We propose that such logical combinations correspond more closely to natural concepts and are more transparent than linear combinations. PSI employs boosting combined with rule mining to encourage learning of non-overlapping (or orthogonal) sets of propositional clauses. A similarity metric is introduced so that textual cases can be compared based on similarity between extracted logical clauses. Interpretability of these logical constructs creates new avenues for user interaction and naturally leads to the discovery of knowledge. PSI's feature extraction approach is compared with the popular dimensionality reduction technique Latent Semantic Indexing (LSI), which uses singular value decomposition to extract orthogonal features that are linear combinations of keywords [7]. Case representations that employ PSI's logical expressions are more comprehensible to domain experts and end-users compared to LSI's linear keyword combinations. Ideally we wish to achieve this expressiveness without significant loss in retrieval effectiveness.

We first establish our terminology for feature selection and extraction, before describing how PSI extracts features as logical combinations. We then describe LSI, highlighting the problem of interpretability with linear combinations. Finally we show that PSI's approach achieves comparable retrieval performance yet remains expressive.

2 Feature Selection and Extraction

Consider the hypothetical example in Figure 1 where the task is to weed out Spam from legitimate email related to AI. To assist with future message filtering these messages must be mapped onto a set of cases before they can be reused. We will refer to the set of all labelled documents (cases) as \mathcal{D} . The keyword-vector representation is commonly used to represent a document d by considering the presence or absence of words [17]. Essentially the set of features are the set of words \mathcal{W} (e.g. "conference", "intelligent"). Accordingly a document d is represented as a pair (\mathbf{x}, y) , where $\mathbf{x} = (x_1, \dots, x_{|\mathcal{W}|})$ is

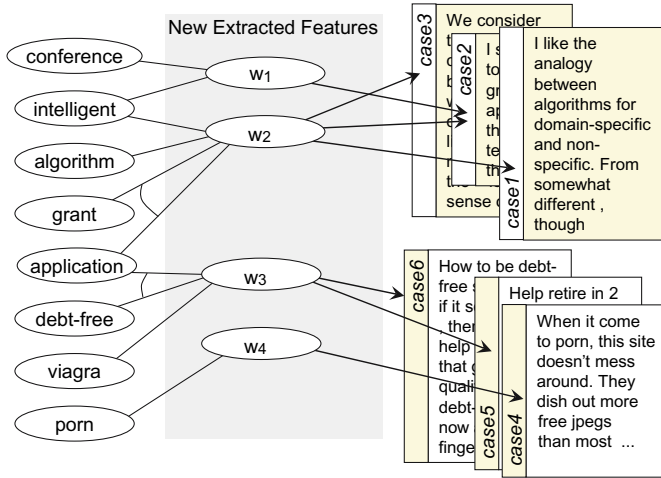


Fig. 1. Features as logical keyword combinations

a binary valued feature vector corresponding to the presence or absence of words in \mathcal{W} ; and y is d 's class label.

Feature selection reduces $|\mathcal{W}|$ to a smaller feature subset size m [20]. Information Gain (IG) is often used for this purpose, where m features with highest IG are retained and the new binary-valued feature vector \mathbf{x}' is formed with the reduced word vocabulary set \mathcal{W}' , where $\mathcal{W}' \subset \mathcal{W}$ and $|\mathcal{W}'| \ll |\mathcal{W}|$. The new representation of document d with \mathcal{W}' is a pair (\mathbf{x}', y) . Selection using IG is the base-line algorithm in this paper and is referred to as BASE. An obvious shortcoming of BASE is that it fails to ensure selection of non-redundant keywords. Ideally we want \mathbf{x}' to contain features that are representative but also orthogonal. A more serious weakness is that BASE's one-to-one feature-word correspondence operates at a lexical level, ignoring underlying semantics.

Figure 1 illustrates a proof tree showing how new features can be extracted to capture keyword relationships using propositional disjunctive normal form clauses (DNF clauses). When keyword relationships are modelled, ambiguities in text can be resolved to some extent. For instance “grant” and “application” capture semantics akin to legitimate messages, while the same keyword “application” in conjunction with “debt-free” suggests Spam messages.

Feature extraction, like selection, also reduces $|\mathcal{W}|$ to a smaller feature subset size m . However unlike selected features, extracted features no longer correspond to presence or absence of single words. Therefore, with extracted features the new representation of document d is (\mathbf{x}'', y) , but $\mathcal{W}'' \not\subset \mathcal{W}$. When extracted features are logical combinations of keywords as in Figure 1, then a new feature $w'' \in \mathcal{W}''$, represents a propositional clause. For example the new feature w''_2 represents the clause: “intelligent” \vee “algorithm” \vee (“grant” \wedge “application”).

3 Propositional Semantic Indexing (PSI)

PSI discovers and captures underlying semantics in the form of propositional clauses. PSI’s approach is two-fold. Firstly, decision stumps are selected by IG and refined by association rule mining, which discovers sets of Horn clause rules. Secondly, a boosting process encourages selection of non-redundant stumps. The PSI feature extraction algorithm and the instantiation of extracted features appear at the end of this section after a description of the main steps.

3.1 Decision Stump Guided Extraction

A decision stump is a one-level decision tree [12]. In PSI, a stump is initially formed using a single keyword, which is selected to maximise IG. An example decision stump formed with “conference” in its decision node appears in Figure 2. It partitions documents into leaves, based on whether or not “conference” appears in them. For instance 70 documents contain the word “conference” and just 5 of these are Spam (i.e. +5). It is not uncommon for documents containing “conference” to still be semantically similar to those not containing it. So documents containing “workshop” without “conference” in the right leaf are still contextually similar to those containing “conference” in the left leaf. A generalised decision node has the desired effect of bringing such semantically related documents closer [19]. Generalisation refines the decision node formed with a single feature w' , to an extracted feature w'' , containing a propositional clause. Typically a refined node results in an improved split (see right stump in Figure 2).

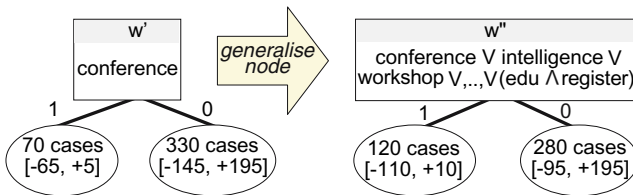


Fig. 2. Node Generalisation

A propositional clause is formed by adding disjuncts to an initial clause containing just the selected feature w' . Each disjunct is a conjunction of one or more keyword co-occurrences with similar contextual meaning to that of w' . An exhaustive search for disjuncts will invariably be impractical. Fortunately the search space can be pruned by using w' as a handle over this space. Instead of generating and evaluating all disjuncts, we generate propositional Horn clause rules that conclude w' given other logical keyword combinations.

3.2 Growing Clauses from Rules

Examples of five association rules concluding in “conference” (i.e. w') appear in Figure 3. These rules are of the form $H \leftarrow B$, where the rule body B is a conjunction of

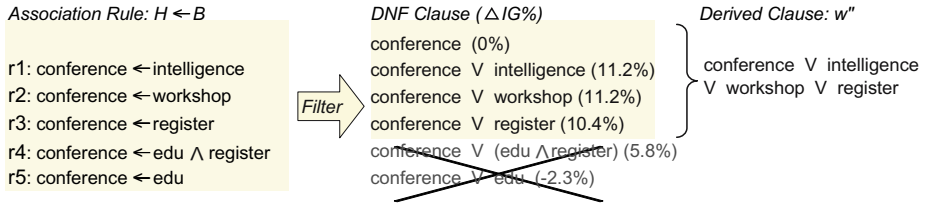


Fig. 3. Growing clauses from selected rules

keywords, and the rule head H is a single keyword. These conjunctions are keyword combinations that have been found to co-occur with the head keyword. Rule bodies are a good source of disjuncts with which to grow our DNF clause w'' , which initially contains only the selected keyword “conference”. However, an informed selection strategy is necessary to identify those disjuncts that are good descriptors of underlying semantics.

The contribution of each disjunct to clause growth is measured by comparing IG of w'' with and without the disjunct (rule body) included in the DNF clause. Disjuncts that fail to improve IG are filtered out by the *gain-filter*. Those remaining are passed onto *gen-filter* where any specialised forms of a disjunct with a lower IG compared to any one of its generalised forms are filtered out. The DNF clauses in Figure 3 show how each rule is converted into a potential DNF clause (difference in IG, used for filtering appear in brackets). The final DNF clause derived once the filtering step is completed is: “conference” \vee “intelligence” \vee “workshop” \vee “register”. We use the Apriori [1] association rule learner to generate feature extraction rules that conclude a selected w' . Apriori typically generates many rules, but the filters are able to identify useful rules.

3.3 Feature Extraction with Boosting

PSI’s iterative approach to feature extraction employs boosted decision stumps (see Figure 4). The number of features to be extracted is determined by *vocabulary_size*. The general idea of boosting is to iteratively generate several (weak) learners, with each learner biased by the training error in the previous iteration [10]. This bias is expressed by modifying weights associated with documents. When boosted stumps are used for feature selection the new document distribution discourages selection of a redundant feature given the previously selected feature [6]. Here, with extracted features, unlike with single keyword-based features, we need to discourage discovery of an overlapping clause given the previously discovered clause. We achieve this by updating document weights in PSI according to the error of the decision stump created with the new extracted feature, w'' , instead of w' .

3.4 Feature Instantiation

Once PSI has extracted new features, textual cases are mapped to a new representation. For a new feature w''_i , let $S_i = \bigvee_j s_{ij}$, be its propositional clause, where $s_{ij} = \bigwedge_k x_{ijk}$

```

 $\mathcal{W}'' = \emptyset; n = |\mathcal{D}|; \text{vocabulary\_size} = m$ 
Algorithm: PSI
Repeat
  initialise document weights to  $1/n$ 
   $w_j = \text{feature with highest IG}$ 
   $\mathcal{W} = \mathcal{W} \setminus w_j$ 
   $w_j'' = \text{GROWCLAUSE}(w_j, \mathcal{W})$ 
   $\mathcal{W}'' = \mathcal{W}'' \cup w_j''$ 
  stump = CREATESTUMP( $w_j''$ )
   $err = \text{error}(\text{stump})$ 
  update document weights using  $err$ 
Until ( $|\mathcal{W}''| = \text{vocabulary\_size}$ )
Return  $\mathcal{W}''$ 

```

Fig. 4. Feature Extraction with PSI

is the j th conjunction in this clause. The new representation of document $d = (\mathbf{x}'', y)$ is obtained by:

$$x_i'' = \sum_j \text{gain_inc}(s_{ij}) * \text{infer}(s_{ij})$$

here `gain_inc` returns the increase in gain achieved by s_{ij} when growing \mathcal{S}_i . Whether or not s_{ij} can be inferred (satisfied) from a document's initial representation $d = (\mathbf{x}, y)$ (i.e. using all features in \mathcal{W}) is determined by:

$$\text{infer}(s_{ij}) = \begin{cases} 1 & \text{if } (\bigwedge_k x_{ijk}) = \text{True} \\ 0 & \text{otherwise} \end{cases}$$

The PSI-derived representation enables case comparison at a semantic (or conceptual) level, because instantiated features now capture the degree to which each clause is satisfied by documents. In other words, satisfaction of the same disjunct will contribute more towards similarity than satisfaction of different disjuncts in the same clause.

4 Latent Semantic Indexing (LSI)

LSI is an established method of feature extraction and dimension reduction. The matrix whose columns are the document vectors $\mathbf{x}_1, \dots, \mathbf{x}_{|D|}$, known as the term-document matrix, constitutes a vector space representation of the document collection. In LSI, the term-document matrix is subjected to singular value decomposition (SVD). The SVD extracts an orthogonal basis for this space, consisting of new features that are linear combinations of the original features (keywords). Crucially, these new features are ranked according to their importance. It is assumed that the m highest-ranked features contain the true semantic structure of the document collection and the remaining features, which are considered to be noise, are discarded. Any value of m less than the

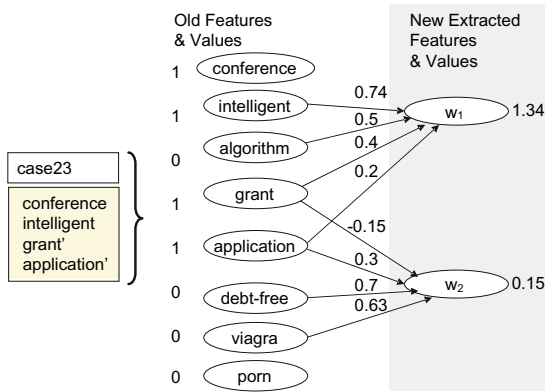


Fig. 5. Feature extraction using LSI

rank of the term-document matrix can be used, but good values will be much smaller than this rank and, hence, much smaller than either $|W|$ or $|D|$. These features form a new lower-dimensional representation which has frequently been found to improve performance in information retrieval and classification tasks [18,22]. Full technical details can be found in the original paper [7].

Figure 5 shows a hypothetical example from the AI-Spam domain (cf. Figure 1). The first extracted feature is a combination of “intelligent”, “algorithm”, “grant” and “application”. Any document containing most of these is likely to be legitimate, so high values of this feature indicate non-Spam. The second feature has positive weights for “application”, “debt-free” and “viagra” and a negative weight for “grant”. A high value for this feature is likely to indicate Spam. The new features are orthogonal, despite having two keywords in common. The first feature has positive weights for both “grant” and “application”, whereas the second has a negative weight for “grant”. This shows how the modifying effect of “grant” on “application” might manifest itself in a LSI-derived representation. With a high score for the first extracted feature and a low score for the second, the incoming test case is likely to be classified as legitimate email.

LSI extracted features are linear combinations of typically very large numbers of keywords. In practice this can be in the order of hundreds/thousands of keywords, unlike in our illustrative example involving just 8 keywords. Consequently, it is difficult to interpret these features in a meaningful way. In contrast, a feature extracted by PSI combines far fewer keywords and its logical description of underlying semantics is easier to interpret. A further difference is that, although both PSI and LSI exploit word-word co-occurrences to discover and preserve underlying semantics, PSI also draws on word-class co-occurrences while LSI does not naturally exploit this information.

5 Evaluation

The goodness of case representations derived by BASE, LSI and PSI in terms of retrieval performance is compared on a retrieve-only CBR system, where the weighted majority

vote from the 3 best matching cases are re-used to classify the test case. A modified case similarity metric is used so that similarity due to absence of words (or words in linear combinations or in clauses) is treated as less important compared to their presence [19].

Experiments were conducted on 6 datasets; 4 involving email routing tasks and 2 involving Spam filtering. Various groups from the 20Newsgroups corpus of 20 Usenet groups [13], with 1000 postings (of discussions, queries, comments etc.) per group, form the routing datasets: SCIENCE (4 science related groups); REC (4 recreation related groups); HW (2 hardware problem discussion groups, one on Mac, the other on PC); and RELPOL (2 groups, one concerning religion, the other politics in the middle-east). Of the 2 Spam filtering datasets: USREMAIL [8] contains 1000 personal emails of which 50% are Spam; and LINGSPAM [16] contains 2893 messages from a linguistics mailing list of which 27% are Spam.

Equal-sized disjoint train-test splits were formed. Each split contains 20% of the dataset and also preserves the class distribution of the original corpus. All text was pre-processed by removing stop words (common words) and punctuation. Remaining words were stemmed to form \mathcal{W} , where $|\mathcal{W}|$ varies from approximately 1,000 in USREMAIL to 20,000 in LINGSPAM. Generally, with both routing and filtering tasks, the overall aim is to assign incoming messages into appropriate groups. Hence, test set accuracy was chosen as the primary measure of the effectiveness of the case representation as a facilitator of case comparison. For each test corpus and each method, the accuracy (averaged over 15 trials) was computed for representations with 20, 40, 60, 80, 100 and 120 features.

Paired t-tests were used to find improvements by LSI and PSI compared to BASE (one-tailed test) and differences between LSI and PSI (two-tailed test), both at the 95% significance level. Precision¹ is an important measure when comparing Spam filters, because it penalises error due to false positives (Legitimate \rightarrow Spam). Hence, for the Spam filtering datasets, precision was tested as well as accuracy.

5.1 Results

Accuracy results in Figure 6 shows that BASE performs poorly with only 20 features, but gets closer to the superior PSI when more features are added. PSI's performance is normally good with 20 features and is robust to the feature subset size compared to both BASE and LSI. LSI clearly performs better for smaller sizes. This motivated an investigation of LSI with fewer than 20 features. We found that 10-feature LSI consistently outperforms 20-feature LSI and is close to optimal. Consequently, 10-feature LSI was used for the significance testing, in order to give a more realistic comparison with the other methods.

Table 5.1 compares performance of BASE and PSI (both 20 features) and LSI (10 features). Where LSI or PSI are significantly better than BASE, the results are in bold. Where LSI and PSI are significantly different, the better result is starred. It can be seen that LSI is significantly better than BASE on 6 of 8 measures and to PSI on 3. PSI is better than BASE on 7 measures and better than LSI on 2. We conclude that the 20-dimensional representations extracted by PSI have comparable effectiveness to the 10-dimensional representations extracted by LSI. Generally, BASE needs a much larger

¹ Precision = TP/(TP+FP) where TP is no. of true positives and FP is no. of false positives.

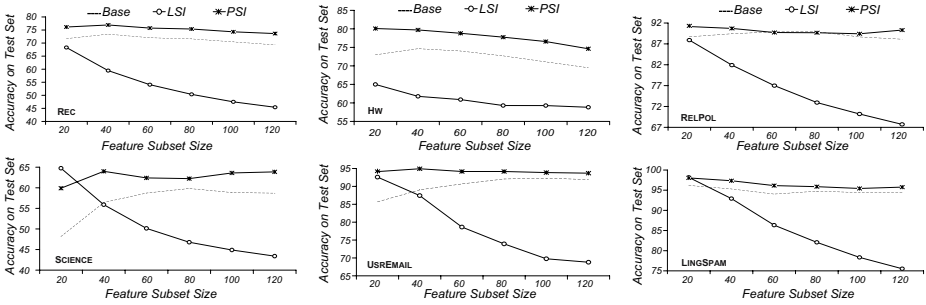


Fig. 6. Accuracy results for datasets

Table 1. Summary of significance testing for feature subset size 20 (10 for LSI)

Algo.	Routing: Accuracy				Filtering: Accuracy (Precision)	
	REC	HW	RELPOL	SCIENCE	USREMAIL	LINGSPAM
BASE	71.7	73.0	88.7	48.1	85.7 (89.5)	94.2 (92.0)
LSI	*78.7	65.5	90.4	*71.8	93.9 (*96.8)	96.8 (89.0)
PSI	76.2	*80.1	91.2	59.9	94.1 (95.2)	95.8 (*92.1)

indexing vocabulary to achieve comparable performance. PSI works well with a small vocabulary of features, which are more expressive than LSI’s linear combinations.

5.2 Interpretability

Figure 7 provides a high-level view of sample features extracted by PSI in the form of logical combinations (for 3 of the datasets). It is interesting to compare the differences in extracted combinations (edges), the contribution of keywords (ovals) to different extracted features (boxes) and the number of keywords used to form conjunctions (usually not more than 3). We see a mass of interconnected nodes with the HW dataset on which PSI’s performance was far superior to that of LSI. Closer examination of this data set shows that there are many keywords that are polysemous given the two classes. For instance “drive” is applicable both to Macs and PCs but combined with “v1b” indicates PC while with “syquest” indicates Mac. Unlike HW, the multi-class SCIENCE dataset contains several disjoint graphs each relating to a class, suggesting that these concepts are easily separable. Accuracy results show LSI to be a clear winner on SCIENCE. This further supports our observation that LSI operates best only in the absence of class-specific polysemous relationships. Finally, features extracted from LINGSPAM in figure 7 show that the majority of new features are single keywords rather than logical combinations. This explains BASE’s good performance on LINGSPAM.

An obvious advantage of interpretability is knowledge discovery. Consider the SCIENCE tree, here, without the extracted clauses indicating that “msg”, “food” and “chinese” are linked through “diet”, one would not understand the meaning in context of a

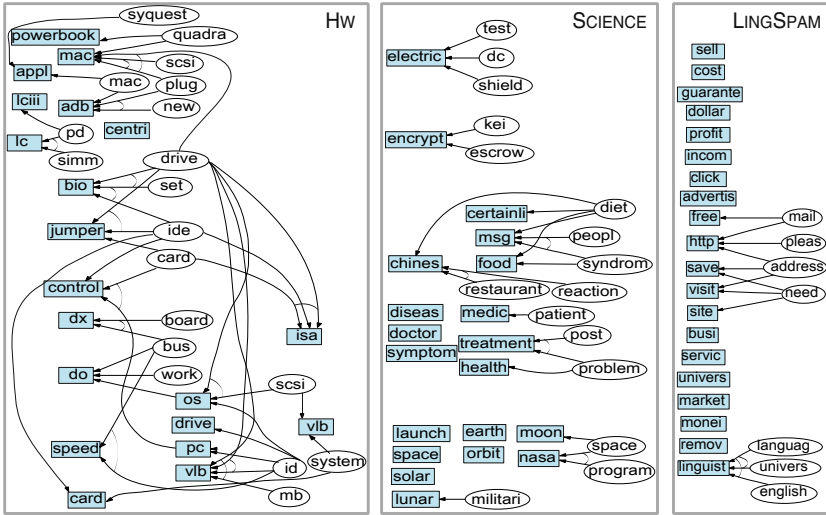


Fig. 7. Logical combinations extracted from datasets

term such as “msg”. Such proof trees, which are automatically generated by PSI, highlight relations between keywords; this knowledge can further aid with glossary generation, often a demanding manual task (e.g. FALLQ [14]). Case authoring is a further task that can benefit from PSI generated trees. For example disjoint graphs involving “electric” and “encrypt” with many fewer keyword associations may suggest the need for case creation or discovery in that area. From a retrieval standpoint, PSI generated features can be exploited to facilitate query elaboration, within incremental retrieval systems. The main benefit to the user would be the ability to tailor the expanded query by deactivating disjuncts to suit retrieval needs. Since clauses are granular and can be broken down into semantically rich constituents, retrieval systems can gather statistics of which clauses worked well in the past, based on user interaction; this is difficult over linear combinations of features mined by LSI.

6 Related Work

Feature extraction is an important area of research particularly when dealing with textual data. In Textual-CBR (TCBR) research the SMILE system provides useful insight into how machine learning and statistical techniques can be employed to reason with legal documents [3]. As in PSI, single keywords in decision nodes are augmented with other keywords of similar context. Unlike PSI, these keywords are obtained by looking up a manually created domain-specific thesaurus. Although PSI grows clauses by analysing keyword co-occurrence patterns, they can just as easily be grown using existing domain-specific knowledge. A more recent interest in TCBR involves extraction of features in the form of predicates. The FACIT framework involving semi-automated index vocabulary acquisition addresses this challenge but also highlights the need for

reliance on deep syntactic parsing and the acquisition of a generative lexicon which warrants significant manual intervention [11].

In text classification and text mining research, there is much evidence to show that analysis of keyword relationships and modelling them as rules is a successful strategy for text retrieval. A good example is RIPPER [5], which adopts complex optimisation heuristics to learn propositional clauses for classification. A RIPPER rule is a Horn clause rule that concludes a class. In contrast, PSI's propositional clauses form features that can easily be exploited by CBR systems to enable case comparison at a semantic level. Such comparisons can also be facilitated with the FEATUREMINE [21] algorithm, which also employs association rule mining to create new features based on keyword co-occurrences. FEATUREMINE generates all possible pair-wise keyword co-occurrences converting only those that pass a significance test into new features. What is unique about PSI's approach is that firstly, search for associations is guided by an initial feature selection step, secondly, associations remaining after an informed filtering step are used to grow clauses, and, crucially, boosting is employed to encourage growing of non-overlapping clauses. The main advantage of PSI's approach is that instead of textual case similarity based solely on instantiated feature value comparisons (as in FEATUREMINE), PSI's clauses enable more fine-grained similarity comparisons. Like PSI, WHIRL [4] also integrates rules resulting in a more fine-grained similarity computation over text. However these rules are manually acquired.

The use of automated rule learning in an Information Extraction (IE) setting is demonstrated by TEXTRISE, where mined rules predict text for slots based on information extracted over other slots [15]. The vocabulary is thus limited to template slot fillers. In contrast PSI does not assume knowledge of case structures and is potentially more useful in unconstrained domains.

7 Conclusion

A novel contribution of this paper is the acquisition of an indexing vocabulary in the form of expressive clauses, and a case representation that captures the degree to which each clause is satisfied by documents. The propositional semantic indexing tool, PSI, introduced in the paper, enables text comparison at a semantic, instead of a lexical, level. Experiments show that PSI's retrieval performance is significantly better than that of retrieval over keyword-based representations. Comparison of PSI-derived representations with the popular LSI-derived representations generally shows comparable retrieval performance. However in the presence of class-specific polysemous relationships PSI is the clear winner. These results are very encouraging, because, although features extracted by LSI are rich mathematical descriptors of the underlying semantics in the domain, unlike PSI, they lack interpretability. We note that PSI's reliance on class knowledge inevitably restricts its range of applicability. Accordingly, future research will seek to develop an unsupervised version of PSI.

Acknowledgements

We thank Susan Craw for helpful discussions on this work.

References

1. R. Agrawal, H. Mannila, R. Srikant, H. Toivonen, and A. I. Verkamo. Fast discovery of association rules. In *Advances in KD and DM*, pages 307–327, 1995. AAAI/MIT.
2. D. Aha, editor. *Mixed-Initiatives Workshop at 6th ECCBR*, 2002. Springer.
3. S. Bruninghaus and K. D. Ashley. Bootstrapping case base development with annotated case summaries. In *Proc of the 2nd ICCBR*, pages 59–73, 1999. Springer.
4. W. W. Cohen. Providing database-like access to the web using queries based on textual similarity. In *Proc of the Int Conf on Management of Data*, pages 558–560, 1998.
5. W. W. Cohen and Y. Singer. Context-sensitive learning methods for text categorisation. *ACM Transactions in Information Systems*, 17(2):141–173, 1999.
6. S. Das. Filters, wrappers and a boosting based hybrid for feature selection. In *Proc of the 18th ICML*, pages 74–81, 2001. Morgan Kaufmann.
7. S. C. Deerwester, S. T. Dumais, T. K. Landauer, G. W. Furnas, and R. A. Harshman. Indexing by latent semantic analysis. *Journal of the American Society of Information Science*, 41(6):391–407, 1990.
8. S. J. Delany and P. Cunningham. An analysis of case-base editing in a spam filtering system. In *Proc of the 7th ECCBR*, pages 128–141, 2004. Springer.
9. G. Forman and I. Cohen. Learning with Little: Comparison of Classifiers Given Little Training. In *Proc of the 8th European Conf on PKDD*, pages 161–172, 2004.
10. Y. Freund and R. Schapire. Experiments with a new boosting algorithm. In *Proc of the 13th ICML*, pages 148–156, 1996.
11. K. M. Gupta and D. W. Aha. Towards acquiring case indexing taxonomies from text. In *Proc of the 17th Int FLAIRS Conference*, pages 307–315, 2004. AAAI press.
12. W. Iba and P. Langley. Induction of one-level decision trees. In *Proc of the 9th Int Workshop on Machine Learning*, pages 233–240, 1992.
13. T. Joachims. A probabilistic analysis of the Rocchio algorithm with TFIDF. Technical report, Carnegie Mellon University CMU-CS-96-118, 1996.
14. M. Lenz. Defining knowledge layers for textual CBR. In *Proc of the 4th European Workshop on CBR*, pages 298–309, 1998. Springer.
15. U. Y. Nahm and R. J. Mooney. Mining soft-matching rules from textual data. In *Proc of the 17th IJCAI*, pages 979–984, 2001.
16. G. Sakkis, I. Androutopoulos, G. Paliouras, V. Karkaletsis, C. Spyropoulos, and P. Stamatoopoulos. A memory-based approach to anti-spam filtering for mailing lists. *Information Retrieval*, 6:49–73, 2003.
17. G. Salton and M. J. McGill. *An introduction to modern IR*. 1983, McGraw-Hill.
18. F. Sebastiani. ML in automated text categorisation. *ACM Computing surveys*, 34:1–47, 2002.
19. N. Wiratunga, I. Koychev, and S. Massie. Feature selection and generalisation for textual retrieval. In *Proc of the 7th ECCBR*, pages 806–820, 2004. Springer.
20. Y. Yang and J. O. Pedersen. A comparative study on feature selection in text categorisation. In *Proc of the 14th ICML*, pages 412–420, 1997. Springer.
21. S. Zelikovitz. Mining for features to improve classification. In *Proc of Machine Learning, Models, Technologies and Applications*, 2003.
22. S. Zelikovitz and H. Hirsh. Using LSI for text classification in the presence of background text. In *Proc of the 10th Int Conf on Information and KM*, 2001.

A Quantitative Comparison of the Subgraph Miners MoFa, gSpan, FFSM, and Gaston

Maarc W., et al., Thorsten Meinl, Igid Fische, and Michael Philippsen

University of Erlangen-Nuremberg, Computer Science Department 2,
Martensstr. 3, 91058 Erlangen, Germany
simawoer@stud.informatik.uni-erlangen.de
{meinl, idfische, philippsen}@cs.fau.de

Abstract. Several new miners for frequent subgraphs have been published recently. Whereas new approaches are presented in detail, the quantitative evaluations are often of limited value: only the performance on a small set of graph databases is discussed and the new algorithm is often only compared to a single competitor based on an executable. It remains unclear, how the algorithms work on bigger/other graph databases and which of their distinctive features is best suited for which database. We have re-implemented the subgraph miners MoFa, gSpan, FFSM, and Gaston within a common code base and with the same level of programming expertise and optimization effort. This paper presents the results of a comparative benchmarking that ran the algorithms on a comprehensive set of graph databases.

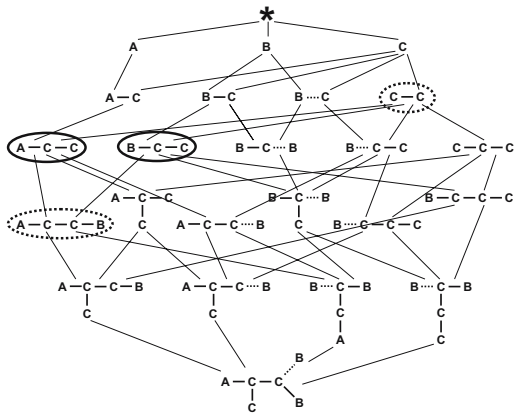
1 Introduction

Mining frequent subgraphs in graph databases is a challenging, especially in the case of large and/or dynamic data sets. The frequent subgraph mining problem is NP-complete [3], and the problem of finding all frequent subgraphs is NP-complete [1].

The main frequent subgraph mining algorithms are based on the idea of generating all frequent subgraphs by adding edges and nodes to a seed frequent subgraph. For each generated frequent subgraph, the algorithm checks if it is frequent. Since the frequent subgraph mining problem is NP-complete, the algorithms are often slow. The main frequent subgraph mining algorithms are based on the idea of generating all frequent subgraphs by adding edges and nodes to a seed frequent subgraph. For each generated frequent subgraph, the algorithm checks if it is frequent. Since the frequent subgraph mining problem is NP-complete, the algorithms are often slow.

A frequent subgraph is a graph that is contained in a graph database. A frequent subgraph is a graph that is contained in a graph database. A frequent subgraph is a graph that is contained in a graph database.

...ed a 1 fe e fag e ...ice he r e e e ... 1 a ea e e ... e
 ...ae ¹E cie fag e ...ie ha e ... e hee ai ... b ... be ...



* is the empty fragment. Each graph is subgraph of all its descendants in the lattice. Subgraphs on one level have the same number of edges.

The dashed C-C-fragment is the common core of the two circled fragments. The new subgraph A-C-C-B can be generated by taking this core and adding the two edges A- and B- that only appear in one of the subgraphs.

Fig. 1. The complete subgraph lattice of the graph shown at the bottom

(A) Purposive refinement. Mining general features from a set of graphs is often a challenging task. The basic approach is to generate all possible subgraphs and then filter them based on some criteria (Fig. 1). A high number of subgraphs is often generated, which can be reduced by using a heuristic (i.e., heuristic) to filter out subgraphs that are not interesting. One of the heuristics used is the *frequency antimonotone principle* (i.e., the frequency antimonotone principle) [4,5]. The idea is to filter out subgraphs that are not frequent enough to be considered interesting.

(B) Efficient enumeration. Generating all subgraphs of a graph is a computationally expensive task. One of the main reasons for this is the large number of subgraphs that need to be generated. However, there are several techniques that can be used to reduce the number of subgraphs that need to be generated. One of these techniques is the *efficient enumeration* technique [6]. This technique uses a heuristic to filter out subgraphs that are not interesting, which reduces the number of subgraphs that need to be generated.

(C) Focused isomorphism testing. Knowledge discovery in graphs is a challenging task. One of the main reasons for this is the large number of subgraphs that need to be generated. However, there are several techniques that can be used to reduce the number of subgraphs that need to be generated. One of these techniques is the *focused isomorphism testing* technique [7]. This technique uses a heuristic to filter out subgraphs that are not interesting, which reduces the number of subgraphs that need to be generated.

Each fragment is a graph. For example, the fragment A-C-C-B is a graph with 4 nodes and 5 edges. The fragment A-C-C-B is a graph with 4 nodes and 5 edges. The fragment A-C-C-B is a graph with 4 nodes and 5 edges.

¹ Similar to the *frequency antimonotone principle* in frequent itemset mining [4,5].

(df) a ... ache ... eed ... e ... ea ... a ce ... beca ... he ... be ... f ... ha ... ha e ... be ... ed ... e ... i ... i ... a ... he ... f he ... a ice (i.e. he ... e f he bigge ... ga h) ... he ea ... i ... i ... a ... i ... (i.e. he ... a ... a ... be ... f ... bg a h ... e ... e e) ... i ... be ad h ... ea che ... The df -a g ... i h ... M Fa [9], gS a ... [10], FFSM [11], a d Ga ... [12] a ac ... he ... b ... be ... (A C) ... i ... e di ... e e ... B ... i ... ce ... i ... i ... di ... c ... e ha a ... i ... i ... be ... e ha a ... he, ... he a h ... a ... e ec a fe da aba e a d ... e e ... be ch a ... ha de ... a e ha he ... ed ... i ... be - ... e ha a c ... e i ... ba ed ... e ec a be ... Si ce di ... e e a h ... e di ... e e da aba e he e i ... ge e a ic ... e. I ... hich f he ... i ... (A C) e f ... be ... de ... hich c ... di ... T ... a e hi g ... e, ... e i ... e e ec a be f he a g ... i h ... a e a ar a be. He ce, ... ea ... e e ... a e ... e ed b ... e f di ... e e ... g a ... i g a g age a d b e ... i a ... f a ... i g c ... i e ... i ... i a ... ech ... g , e c.

I ... hi ... a e ... e e ... a ... bia ed a d de a red c ... a ... f he f ... f ag e ... i e. M Fa, gS a , FFSM a d Ga We ... e e ed he a f ... c a ch ... i g a c ... g a h f a e ... , i.e. a ... e he a e g a h da a ... c ... e. I ... ec ... 2 ... e b ... ie cha ac e i e h ... he e a g ... i h ... e b ... be ... (A C). Sec ... 3 c ... a ... he ... ar b d ... f hi ... a e : he de a red e ... e i ... e a e a a ... f he f ... c ... e a ...

2 Distinctive Ideas of MoFa, gSpan, FFSM, and Gaston

A f ... f ag e ... i e ... ge e a , ... di ec ed g a h ... i h a be ed ... de a d edge . The a a e e ... ic ed ... di g c ... ec ed ... bg a h a d , a e e he a ice a ... e ... ed be f ... e i ... de h ... de .

MoFa (M ec e F ag e M i e , b B , ge a d Be h d i 2002 [9]) ha bee a ge ed ... a d ... ec a da aba e , b ... i ca a ... be ... ed f ... a , b i , a g a h . M Fa ... e a e beddi g (b ... h ... de a d edge) . E ... e ... i ... e ... ic ed ... h ... e f ag e ... , ha ac a ... a ... ea ... i he da aba e . I ... h i ... e ... i he da aba e ca chea ... be d ... e b ... e i g he he a e beddi g ca be ... ed ... i he a e a . M Fa ... e a f ag e ... - ca ... be i g che e ... ed ce he ... be ... f ... e e ... ge e a ed f ... a f ag e : M Fa c ... he ... de ... f a f ag e ... acc ... di g ... he e e ce ... hich he ha e bee added . Whe a f ag e ... i e e ded a ... de n , a e ... e e e ... a ... c c ... a n ... a ... de bigge ... ha n . M ... e e , a e e ... i ... ha g ... f ... he a e ... de n a e ... de ed acc ... di g ... i c ea i g ... de a d edge a be . A h gh hi ... ca ... de i g he ... , M Fa ... i ge e a e ... a ... huc f ag e ... a d he ... e ... a da d ... h i ... e i g ... ed ... i ca e .

gSpan (g a h -ba ed S b ... c ... e a e ... b Ya a d Ha ... 2002 [10]) ... e a ca ... i ca ... e ... e a ... f ... g a h , ca ed df -c de . A df - , a e a ... f a g a h de ... e a ... de ... i ... hich he edge a e ... i ed . The c ... ca e a ... f edge ... e e a ... i ... i ha ... de ... i he g a h ' df -c de . Re ... e e ... ge e a - ... i ... i ... e ... ic ed b gS a ... a ... : F ... , f ag e ... ca ... be e e ded a ... de ha ... he ... , ... f he df - ... ee . Sec ... d , f ag e ... ge -

... a 1... 1 g ided b... cc... e ce 1... he a... ea, a ce 1... . Si ce he e... 1 g... e ca... f... e e... 1... hic f ag... ge e a 1... , gS a c... e he ca... ica (e ic g a hica... a e) df -c de f, each e... e b... ea... f a... ie f e... a 1... . Re... e... 1 h... - 1 1 a df -c de ca be... ed. Si ce 1... ead f e bedd i g, gS a... e a... ea, a ce 1... f, each f ag... e... e... ic 1... bg a h 1... , h 1... e 1 g... bed... e... a g a h 1... he e a... ea, a ce 1... .

FFSM (Fa Ee e Sbg a h Mi 1 g, b Ha , Wa g, a d P... 2003 [11])... e... e... g a h a... ia ge... a ice (... de abe... he diag... a, edge abe... e... he e). The... a 1 -c de 1... he c... ca e a 1... f a 1... e... ie, ef... igh a d 1 e b... 1 e. Ba ed... e ic g a h ic... de 1 g, 1... , hic g a h ha e he... a e ca... ica c de (CAM Ca... ica Ad ace c Ma... 1). Whe... FFSM... 1... a ice f f ag... e... ge e a e... e... e... ,... a... e... e... c... e... e... . FFSM a... eed a... e ic ed e... e... 1... e a 1... : a... e edge... de a 1... a... be added... f a CAM. Af e... e... e... ge e a 1... , FFSM e... e... a 1... 1 e... chec... he he... a ge e a e d... a 1... 1... ca... ica f... . If... , 1 ca be... ed. FFSM... e... e bedd i g... a id e... ic 1... bg a h 1... , h 1... e 1 g. H... e e... , FFSM... e... he... a ch i g... de, edge a e ig... ed. Thi... he... eed i g... he... 1 a d e... -... 1... e a 1... 1 ce he e bedd i g 1... f... e f ag... e... ca be ca c a ed b... e... e a 1... . he... de... .

Gaston (GA h/Se... e ce/Tee e... ac 1ON, b Ni... e a d K... 2004 [12])... e... e... e bedd i g, ... ge e a e... e... e... ha ac a... a... ea a d... ach e e fa... 1... , h 1... e 1 g. The... a 1 igh 1 ha he e a... e e cie... a... e... e... e a e a h a d (... -c c ic) ... ee. B... c... ide 1 g f ag... e... ha a... e a h... , ee... , a d b... , ceed i g... ge e a g a h... 1 h c ce a... he e d, a a ge f ac 1... f he... ca be d... e e cie... . O... 1 ha a... ha e, Ga... face... he NP-c... e... e... f he... bg a h 1... , h 1... , be... . Ga... de... e a g... ba... de... c ce-c... 1 g edge a d... ge e a e... h... ec ce ha a... e a ge... ha... he a... e. D... ica e de ec 1... 1 d... e 1... ha e : ha h i g... , e... a d a g a h 1... , h 1... e f... a d... ica e de ec 1... .

F... gSPa... a d M Fa... e e a e... e... 1... e 1... ha a... e de c 1 bed 1... ec 1... 3.5.

3 The Comparison

In the following section we compare the performance of the algorithms based on a series of experiments on a 1.3 GHz processor. The results are presented in the following section.

3.1 Setup of Experiments

The experiments were conducted on a 64-bit Linux system because of the high performance of the Linux system. The biggest database was used because of the high performance of the database. The system configuration was: Intel Core 2 Duo E6700, 2 GB RAM, 1.3 GHz processor, 10 GB of RAM. The hardware used was IBM's Java VM.

MachLearn (JVM) 1.4.2 because it is distributed here because it is free of charge. The available hardware is a Pentium III 1.3 GHz. The available software is Sun JVM 1.4.2 and IBM JVM 1.4.2. The available hardware is a Pentium III 1.3 GHz. The available software is Sun JVM 1.4.2 and IBM JVM 1.4.2. The available hardware is a Pentium III 1.3 GHz. The available software is Sun JVM 1.4.2 and IBM JVM 1.4.2.

We chose Java as the programming language because it is the most popular language for building applications. The available hardware is a Pentium III 1.3 GHz. The available software is Sun JVM 1.4.2 and IBM JVM 1.4.2.

Because the available hardware is a Pentium III 1.3 GHz. The available software is Sun JVM 1.4.2 and IBM JVM 1.4.2. The available hardware is a Pentium III 1.3 GHz. The available software is Sun JVM 1.4.2 and IBM JVM 1.4.2.

Dataset	# molecules	average size # edges	largest molecule # edges	# node labels
IC93	1,283	28	81	10
HIV	42,689	27	234	58
NCI	237,771	22	276	78
PTE	337	26	213	66
CAN2DA99	32,557	28	236	69
HIV CA	423	42	196	21
HIV CM	1,083	34	234	27

Fig. 2. The molecular datasets used for testing and their sizes. There are always four edge labels in molecules.

Our choice of the available hardware is a Pentium III 1.3 GHz. The available software is Sun JVM 1.4.2 and IBM JVM 1.4.2. The available hardware is a Pentium III 1.3 GHz. The available software is Sun JVM 1.4.2 and IBM JVM 1.4.2.

² <http://www-128.ibm.com/developerworks/java/jdk/index.html>
³ <http://java.sun.com/>
⁴ <http://www.sgi.com/products/servers/altix/index.html>
⁵ <http://www.bea.com/framework.jsp/content/products/jrocket/>
⁶ http://dtp.nci.nih.gov/docs/aids/aids_data.html
⁷ <http://cactus.nci.nih.gov/ncidb2/download.html>
⁸ See [16] and <http://web.comlab.ox.uk/oucl/research/areas/machlearn/PTE/>. The dataset we used was provided by Siegfried Nijssen.
⁹ http://dtp.nci.nih.gov/docs/cancer/cancer_data.html

CAN2DA99, the e da a e, a e, a he, a c, a ed, he c, e e HIV, the NCI da a e.

3.2 Hotspots

Section 2 has already highlighted the frequent, high-level, but basic (A-C) and high-level, but basic, He ce, e, h, he, i ed, i b, b, e, ce, age f, each a a d each a g, i h, a, ea, e e, ha a, d, e bef, e i, he i e a, e. We used Quest JP, be, a P, e¹⁰, the PC, i h, the SUN JVM f, i, i g a, the IC93 da a e, i ha, i i, f 5%, see Fig. 3. Using a, e, d, he, i e a, e, e, the biggest database, ha a e, a age a b e f, h e, e, i e: IC93 and HIV CA + HIV CM.

	IC93				HIV CA+CM			
	MoFa	gSpan	FFSM	Gaston	MoFa	gSpan	FFSM	Gaston
Duplicate filtering/pruning	11.3%	3.1%	0.1%	1.8%	12.3%	1.4%	0.2%	1.0%
Support computation	9.3%	62.9%	3.7%	87.8%	9.6%	70.7%	3.3%	95.9%
Embedding list calculations	19.1%	-	60.4%		18.1%	-	62.7%	
Extending of subgraphs	29.9%	17.3%	10.2%		31.1%	14.9%	8.1%	
Joining of subgraphs	-	-	0.1%	-	-	-	0.1%	-

Fig. 3. The table shows the main parts of the subgraph mining process and how much time (relative to the total runtime) each of the four algorithms spends for them

Filtering/pruning duplicates a, a, i, e i, he h e, b, g a h, i g, ce, (0.1% - 12.3% f he, a, i e). F, M Fa, he i e c, a i, b, h he g a h, h, e, f, a, e a d g a h, a d he d e i, f e, i, ha d, c, i h he, c, a, i g, e.

Support computation or Embedding list calculation he e he a g, i h, e d, f he, i e. Using e beddi g i (M Fa a d FFSM) e a d, b e, i, c, a i, b, c a i g he i e e, i e. A h gh M Fa? 19.1% f, IC93 ee fa e, ha FFSM? 60.4% f, IC93, b, h a g, i h, ha e, e, ab, he a e, b e, f, e c, d, i, h, a. If e beddi g i, a e d (g S a), e e i e, b g a h, h, e, a, e, e c e a. F, G a, i, i, i, b e, e a a e, i e f, c, a i, e beddi g i, c a i, a d he e e, i, f f a g e. The 87.8% f, IC93 i c d e G a, a b i, e, g e e a e a h a d e e.

Extending or joining subgraphs a e ab, he a e i e i, M Fa, g S a a d FFSM. J, i i g i, d, e b FFSM a d i e, chea c, a e d, he e, e, i, ce.

The, b e, f, HIV CA + CM d, d, e, ch f, he, b e, e a, e d f, IC93.

¹⁰ <http://www.quest.com/jprobe/>

he age da aba e a d Ga ... eed he ... e ... fa ag, i h ... ee
 Fig. 5 ... he ef .

The e i a ... e ... e cea ... i e a i g a ... g he f ... ag, i h ... :
 M Fa i a a ... he ... e . O he big da a e , FFSM i he ec d ... e
 ag, i h ... IC93 i i fa e ha gS a . The e ... f hi IC93 e
 e a he e ... e i [15]. The i e ... ea ... h gS a i ... he
 IC93 da a e i he g ... i g ... be f ... bg a h i ... hi ... e ... gS a ha
 ... d a he e ... a e (... e ha 37,000). A ... he ag, i h ... e
 e beddi g i ... hich eed ... he e ... e ecia f ... a gef ag e ... O he
 a ge da a e , h ... e e , gS a i fa e ha FFSM. Ga ... i he fa e f
 a ag, i h ... e ce a ... e ... a e ... he c ... e e NCI da a e . A
 ... f ... hi ... a be he a ... f b ... ee i g beca e f he a ge ... be
 f e beddi g . A ... a ... d ... beca e f ... e f e e ga bage c ec i ...
 a be he ca e . The f ag e ... f d a e h ... e e , a he ... a ... ha gS a
 ge b ... i h chea e ...

The ... be f f d d i ca e (igh c ...) g i e a i igh i ... he
 ... e f he f ag e ... e e e ... echa i ... A ... di e e ... i g ech i e
 ... i i i e he ... be f d i ca e , b a h ... i ec i . 3.2 he a e ... a
 ... e a . Ga ... i a i d e ... d ce d i ca e f ... -c c i c g a h .
 O he ... he ha d FFSM' a d M Fa' e e i ... e h d a d ... i g , e
 ... ee ... be he ea e . F , FFSM a ... a ea i e i e e d i ... e i g .
 he e d i ca e (ee a b e 3) , hich i ... 0.1% i e Ga ... , i dica e ha he
 ca ... i ca , e , e e a i ... i e , e c i e .

Ne he e ... c ... i a a i g ... a e a ec, ded ba ed
 ... he SUN JVM. We f e e ... ca ed he ga bage c ec ... a d ec, ded he
 ... a i ... he a i e . Thi d e ... e ce a i g i e he e ac a e f he e -
 ... c ... i ... , b i a e g d a ... i a i . Beca e i ... d ... he
 ... i e d a a i ca ... he a e f ... he HIV da a e ... e e ec, ded. A ca
 be ee i Fig. 5, gS a eed he ea ... e ... a i d e ... ee beddi g
 i ... A h gh M Fa ... e b hedge a d ... de i he e beddi g i ... he ea
 FFSM ... e he ... de , M Fa ... i eed e ... e Thi i beca e M Fa
 ... eed ... e i each ... de f he ea ch ee he e beddi g f ... e b-
 g a h , hi e i FFSM a ea ch ee ... de c ... i ... f a ... bg a h ... ge he
 i h he i e beddi g . Ga ... eed he ... e ... , beca e e beddi g i
 f ... a e f ag e a e b i ba ed ... he e beddi g i ... f he a e . E -
 e ... i ... he a e ? e beddi g i a e ... ed i h he chi d e . The ef , e ,
 he i e f he e beddi g i d e a ... de e d ... he ... be f chi d e a
 f ag e ha . Thi e ... i he i e f he c ... e f ... , a e .

Fi a ... he ca ab i i f he ag, i h ... f , i c ea i g da aba e i e a
 e ed (Bea JVM, A i) , ee Fig. 5, igh . The c ... e e NCI da aba e a
 ... i i ... 119 i ce e f 2,000 , a d ... e ec ed ... ec e . F , 5% ... e
 ha e e ed he e f ... a ce f ... a i ... be f he NCI da aba e , each . b e
 c ... i g f ag . i g ... be f he e i ce . A ... b i ... c ... i ... i , ha
 a ag, i h ... ca e i ea ... i h he da aba e i e , b ... i h di e e fac
 The ... , i i g e ... i , ha i hi e Ga ... i a a ... e ha gS a

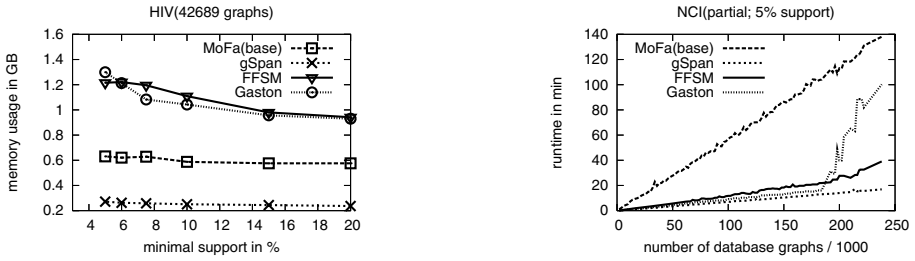


Fig. 5. Memory usage on the HIV database and the runtime in dependence of the database size on the complete NCI database

which are not the case in a real life setting. We have experimented with the complete NCI database and the artificial graph database. The results show that the complete NCI database is not suitable for the complete NCI database. The artificial graph database is suitable for the complete NCI database. The results show that the complete NCI database is not suitable for the complete NCI database. The artificial graph database is suitable for the complete NCI database.

3.4 Tests on Artificial Graph Databases

Real world data are often very large and complex. For example, the complete NCI database is a very large and complex database. The artificial graph database is a very large and complex database. The artificial graph database is a very large and complex database. The artificial graph database is a very large and complex database. The artificial graph database is a very large and complex database.

Next, we have tested the performance of the algorithms on the artificial graph database. The results show that the complete NCI database is not suitable for the complete NCI database. The artificial graph database is suitable for the complete NCI database. The results show that the complete NCI database is not suitable for the complete NCI database. The artificial graph database is suitable for the complete NCI database. The results show that the complete NCI database is not suitable for the complete NCI database. The artificial graph database is suitable for the complete NCI database.

- C...a...c... beief e beddi g i... d... c...ide ab...eed...
he ea ch f...f e e f a g e... E e h gh gS a d e... e he ,
i i c... e i i e... Ga... a d FFSM. O... if he f a g e... be... e a g e
(i e i... he IC93 da a e), gS a fa... O... he... he ha d, e beddi g
i... ca ca e... be... if... e... gh... e... i a a a b e... if he... e...
h... gh... i... high e... gh.
- The... e... f... he... i g... a egie... a i d d... i ca e i... he...
i... a... fa... The ge... e a i... f ca dida e a d... /e beddi g i...
c... a i... a e... ch... e c i i ca.
- U i g ca... i ca... e... e a i... f... de ec i g d... i ca e i... e e cie... ha
d i g e... i c i g a h i... h i... e... E e be... e i... he c... e e a... i da ce
f d... i ca e f a g e... ge... e a i... i e Ga... d e (a... ea f... -c c i c
f a g e...).
- A a g... i h... ca e i ea... i h... he da aba e... i e h... gh... i h d i e e
fa... .
- De e d i g... he... ed Ja a V i... a Mach i... e... ca... e i e d i e...
Thi... be... ca... be... ed b... he a g... i h... he... ef.
- P... e... e f... a ce i... e... e... h i g... A h... gh M Fa i... he... e a g...
i h... i a... e... i... e... ch... e f... c i... a i... ha... he... he... i e... f...
... ec... a... da aba e a d b i... che... i ca... e i... .

I i... e cea, he e he de e... e... f f e... e... bg a h... i i g... i e ad
i... he f... e... P... i b e d i e c i... a... e d i... b... e d... a a e... ea ch... e... e... e
... e... a d e... f... a ce i... . E... i g... e... a... i ca... i... a... ea i e... ec ed...
ead... e... i... gh... .

References

1. Fischer, I., Meinel, T.: Subgraph Mining. In Wang, J., ed.: Encyclopedia of Data Warehousing and Mining. Idea Group Reference, Hershey, PA, USA (2005)
2. Washio, T., Motoda, H.: State of the Art of Graph-based Data Mining. SIGKDD Explorations Newsletter **5** (2003) 59–68
3. McKay, B.: Practical graph isomorphism. Congressus Numerantium **30** (1981)
4. Agrawal, R., Imielinski, T., Swami, A.N.: Mining Association Rules between Sets of Items in Large Databases. In Buneman, P., Jajodia, S., eds.: Proc. 1993 ACM SIGMOD Int’l Conf. on Management of Data, Washington, D.C., USA, ACM Press (1993) 207–216
5. Zaki, M.J., Parthasarathy, S., Ogihara, M., Li, W.: New Algorithms for Fast Discovery of Association Rules. In Heckerman, D., Mannila, H., Pregibon, D., Uthurusamy, R., Park, M., eds.: In 3rd Int’l Conf. on Knowledge Discovery and Data Mining, AAAI Press (1997) 283–296
6. Cook, D.J., Holder, L.B.: Substructure Discovery Using Minimum Description Length and Background Knowledge. J. of Artificial Intelligence Research **1** (1994) 231–255
7. Inokuchi, A., Washio, T., Motoda, H.: An apriori-based algorithm for mining frequent substructures from graph data. In: PKDD ’00: Proceedings of the 4th European Conference on Principles of Data Mining and Knowledge Discovery, London, UK, Springer (2000) 13–23

8. Kuramochi, M., Karypis, G.: Frequent subgraph discovery. In: Proceedings of the IEEE Intl. Conf. on Data Mining ICDM, Piscataway, NJ, USA, IEEE Press (2001) 313–320
9. Borgelt, C., Berthold, M.R.: Mining Molecular Fragments: Finding Relevant Substructures of Molecules. In: Proc. IEEE Int'l Conf. on Data Mining ICDM, Maebashi City, Japan (2002) 51–58
10. Yan, X., Han, J.: gSpan: Graph-Based Substructure Pattern Mining. In: Proc. IEEE Int'l Conf. on Data Mining ICDM, Maebashi City, Japan (2002) 721–723
11. Huan, J., Wang, W., Prins, J.: Efficient mining of frequent subgraphs in the presence of isomorphism. In: Proceedings of the 3rd IEEE Intl. Conf. on Data Mining ICDM, Piscataway, NJ, USA, IEEE Press (2003) 549–552
12. Nijssen, S., Kok, J.N.: Frequent Graph Mining and its Application to Molecular Databases. In Thissen, W., Wieringa, P., Pantic, M., Ludema, M., eds.: Proc. of the 2004 IEEE Conf. on Systems, Man and Cybernetics, SMC 2004, Den Haag, The Netherlands (2004) 4571 – 4577
13. Institute of Scientific Information, Inc. (ISI): Index chemicus - subset from 1993 (1993)
14. Nijssen, S., Kok, J.N.: A quickstart in frequent structure mining can make a difference. Technical report, Leiden Institute of Advanced Computer Science, Leiden University (2004)
15. Huan, J., Wang, W., Prins, J.: Efficient mining of frequent subgraphs in the presence of isomorphism. Technical report, Department of Computer Science at the University of North Carolina, Chapel Hill (2003)
16. Srinivasan, A., King, R.D., Muggleton, S.H., Sternberg, M.: The predictive toxicology evaluation challenge. In: Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence (IJCAI-97). Morgan-Kaufmann (1997) 1–6
17. Yan, X., Han, J.: Closegraph: Mining Closed Frequent Graph Patterns. In: Proc. of the 9th ACM SIGKDD Int'l Conf. on Knowledge Discovery and Data Mining, Washington, DC, USA, ACM Press (2003) 286–295
18. Meinel, T., Borgelt, C., Berthold, M.R.: Discriminative Closed Fragment Mining and Perfect Extensions in MoFa. In Onaindia, E., Staab, S., eds.: STAIRS 2004 - Proc. of the Second Starting AI Researchers' Symp. Volume 109 of Frontiers in Artificial Intelligence and Applications., Valencia, Spain, IOS Press (2004) 3–14
19. Hofer, H., Borgelt, C., Berthold, M.R.: Large Scale Mining of Molecular Fragments with Wildcards. In: Advances in Intelligent Data Analysis. Number 2810 in Lecture Notes in Computer Science, Springer (2003) 380–389
20. Meinel, T., Borgelt, C., Berthold, M.R.: Mining Fragments with Fuzzy Chains in Molecular Databases. In Kok, J.N., Washio, T., eds.: Proc. of the Workshop W7 on Mining Graphs, Trees and Sequences (MGTS '04), Pisa, Italy (2004) 49–60

Efficient Classification from Multiple Heterogeneous Databases*

Xiaojin Yi and Jianfei Han

University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA
{xyin1, hanj}@uiuc.edu

Abstract. With the fast expansion of computer networks, it is inevitable to study data mining on heterogeneous databases. In this paper we propose *MDBM*, an accurate and efficient approach for classification on multiple heterogeneous databases. We propose a regression-based method for predicting the usefulness of inter-database links that serve as bridges for information transfer, because such links are automatically detected and may or may not be useful or even valid. Because of the high cost of inter-database communication, *MDBM* employs a new strategy for cross-database classification, which finds and performs actions with high benefit-to-cost ratios. The experiments show that *MDBM* achieves high accuracy in cross-database classification, with much higher efficiency than previous approaches.

1 Introduction

The rapid growth of heterogeneous databases, especially the emergence of geospatial data, has led to a wide range of applications, such as location-based services, network-based services, and data mining. For example, bibliographic classification, recommendation systems, and data mining are all important applications. Data integration is a key challenge [5,11] in the field of data integration. However, effective integration of heterogeneous databases is a challenging task. In the past, data integration [3,7,8,12,13] and data mining on heterogeneous databases have been studied. The first step is to identify the data sources and their relationships. In this paper, we propose a new method for identifying useful inter-database links. The main idea is to use a regression-based method to predict the usefulness of inter-database links. The experiments show that our method achieves high accuracy in cross-database classification, with much higher efficiency than previous approaches.

* The work was supported in part by the U.S. National Science Foundation NSF IIS-02-09199/IIS-03-08215. Any opinions, findings, and conclusions or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the funding agencies.

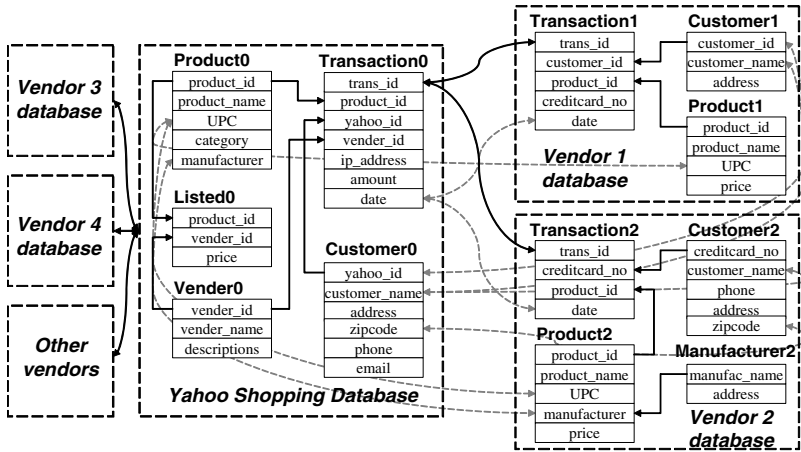


Fig. 1. Databases of Yahoo shopping and vendors

... da aba e f d i e e e d ... I h i e a e e h e , ... e a i i .
 he Yah . h i g da aba e i ca e d , h e e e a e . . .
 The g a f c a -da aba e c a i ca i i b i da acc a e c a i e
 f , e d i g h e c a i a b e f a g e e e .

The e a e . . . a l , c h a e g e i c a -da aba e c a i ca i . The . . .
 i he **data heterogeneity problem**. T a f e i f , a i a c a h e e .
 g e e . da aba e , e e . . . d e c , h i c h a e i . b e
 e e . a c h e d a i b e (a i Fig e 1) a d c a e e a b d g e f i f ,
 a i i , a f e . The e a e . a . d i e . . h i i e , c h a c h e a . a i g
 [5] a d . i i g da aba e . c e [4]. H e e , . . e i . d e c e d . a b e
 a g e a d . e i e c e c . e a e d . b e c . F e a e e ,
 → a c c e c d i e e e i h a e a e , a d , . .
 → a e a d . a e . i e . b e f i e d
 e e . The e c d c h a e g e i he **efficiency problem**. I i f e e e i e
 . . a f e i f , a i i . b e e e . da aba e , h i c h a b e f a f . e a c h
 . h e h i c a . T h e e . . b e a b e . b i d a c c a e c a -da aba e c a i
 e . i h a . i e -da aba e c . . i c a i c . a . . i b e . I h i a e
 e . . e (M i -Da aba e M i e) , a e c i e a d a c c a e a . a c h
 f , c a i ca i . a c i e h e e . g e e . da aba e .

The . . . c a i b i . f h i a e i . . . e a a . a c h f , e d i g
 i g h e f e f e f i . A e i e d a b e , e i e c a e a d . e f
 f e a e , h i e e e h e . a b e e e a d . . a d b . d e . . h e c a i
 c a i . c e d e . We d e e h e f e f a i a h e a i . i f , a i i .
 g a i f a . f e a e g e e a e d b . . a g a i g i f , a i i . h i g h h i i . We
 . . e a e g e i . -b a e d a . a c h f , b i d i g a . d e . . e d i g h e f e f e
 . f i . b a e d . . . e i e f i . O e e i e . h . h a h i a . a c h
 a c h i e e a . a b h i g h e d i g h a c c a c .

O e e c d c a i b i . i A . a a . a c h e . .
 e a i . a (. d e .) c a i ca i . [1,9,10,14], M D B M a . . e . e -b a e d
 c a i ca i . A e i e . a . a c h e b i d . e b e a , c h i g f , e d i c a e (.

1 e a) 1 h ghe 1 f a 1 gai (F 1 gai), 1 de b 1 d acc a e
 e . A h gh h 1 a eg 1 e ec 1 e 1 1 g e da aba e , 1 a ead
 high 1 e -DB c 1 ca 1 c 1 1 -da aba e ca 1 ca 1 . Wi h he
 edic 1 de f gai f e f 1 , MDBM ca edic he gai a d c 1 f
 each ac 1 f ea chi g f edica e . The a eg f
 a a e ec he ac 1 1 h ghe gai - -c a 1 , i.e., he ac 1 f e
 ice e 1 f gai . I ca a che e a e a gai 1 h ch e c . O
 e e 1 e h ha MDBM a che e a high acc a e 1 a a che ,
 b 1 ch e e e c 1 b h 1 g 1 e a d i e -DB c 1 ca 1 .
 The e f he a e 1 ga 1 ed a f . We di c e e a ed 1
 Sec 1 2. Sec 1 3 de c 1 be he a a ch f b 1 d i g edic 1 de f
 e f e f 1 . We de c 1 be he a eg f ec 1 ca c 1 -da aba e ca 1 -
 ca 1 1 Sec 1 4. We e e e 1 ca e a a 1 1 Sec 1 5, a d c 1 c de
 h 1 d 1 Sec 1 6.

2 Related Work

The ad 1 a a f 1 1 g 1 e da aba e 1 e g a e he
 da aba e [5,11], he a da a 1 1 g a g 1 h . H e e , 1 1 f e ha d
 1 e g a e he e ge e da aba e 1 g a e e h e da aba e a he
 1 e , beca e f b h e c e c a d 1 ac c e . Th 1 -da aba e
 1 1 g e eed e c e a a che ha ca d ce g d 1 1 g e 1 h
 1 e -da aba e c 1 ca 1 c 1 .

Di b ed da a 1 1 g e c e d ch a e 1 1 he a e e a e a ,
 hich a 1 a d c e 1 g e d e f a da a e ha 1 d i b ed a d i e -
 e 1 e . The e a e e f d i b ed da a : (1) h 1 a a 1 1 ed
 da a , 1 hich da a ab d i e e b e c 1 h a e a 1 b e a e d b
 d i e e 1 e ; (2) e 1 ca a 1 1 ed da a , 1 hich d i e e a 1 b e f he
 a e e f b e c a e d a d i e e 1 e . E i he a f d i b 1 d i d e
 he c f a a b e 1 d i e e a . Di b ed da a 1 1 g a -
 a che f h 1 a a 1 1 ed da a 1 c de e a - e a 1 g [3] ha e ge
 de b 1 f d i e e 1 e , a d 1 ac e e 1 g e ch 1 e 1 c d i g
 de c 1 e e [8] a d a c 1 a 1 e 1 1 g [7]. Th e f e 1 ca a 1 1 ed
 da a 1 c de a c 1 a 1 e 1 1 g [12] a d k - e a c e 1 g [13]. Di b ed
 da a 1 1 g a e - f a ed da a a b e d a d i e e 1 e . I 1
 f da e a d i e e f c -da aba e da a 1 1 g , hich
 1 e he e ge e da aba e , each c a 1 g a e f 1 e c e c ed e a 1 .

The e a e a d i e e a 1 a (e - de) ca 1 ca 1 [1,9,10,14],
 hich a 1 a b 1 d i g acc a e ca 1 e 1 e a 1 a da aba e . S ch a g -
 1 h e a c h a g d i e e e a 1 f e f edica e , b a a f e 1 g
 1 f a 1 a c e a 1 . The e i h e b 1 d e b a d d i g g d i e a (e
 edica e), b 1 d de c 1 e e e c 1 e . S ch a a che ha e e e
 b e e c e a d acc a e 1 1 g e - da aba e c e a 1 . H e e , 1 -da aba e
 ca 1 ca 1 , he a ha e high 1 e - da aba e c 1 ca 1 c 1 , beca e
 he f c d i g gai f 1 e a b h ch da a eed b e

... a fe, ed. MDBM f... he1... a1... h1... h fca1ca1. (e-ba ed, geed ca1ca1.), b ad... a e... a eg ca ed... hich ca achie ea high acc, ac... i h... ch... e, c... .

3 Predicting Usefulness of Links

3.1 Propagating Information Across Databases

In [4] a e cie a... ach1... ed... ide if... iabe a... i b e... i a e a... i a da aba e. I... ai idea1... c... e he e... e b a ce f e... f a e... f d i e e... a... i b e... a d i a chie e g... d ca a b i... i h a a... i g e ch i e. MDBM... e h i a... ach... d a... i a b e a... i b e a c... da aba e. F... a... i b e A_1 a d A_2 i d i e e... da aba e, i f a i g i ca... i... (a e a 25%) f a e... f A_1 a e... i a b e... A_2 ,... h e f A_2 a e... i a b e... A_1 , he MDBM a... e he e i a... i b e... A_1 a d A_2 . Thi a... ach ha... he i... a... i... ha i ca... d e c... i... e... i... A... e c... che a... a ch i g a... ach [5] ha ca... d e c... c... e... i... (e.g., *firstname + lastname* → *name*) ca... a... b e e a... i... i... e g a e d... i... MDBM.

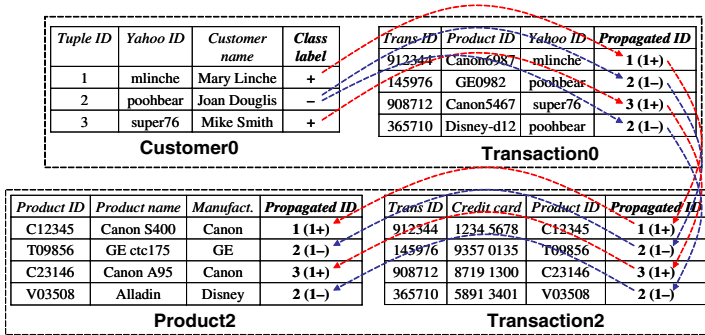


Fig. 2. Example of Tuple ID Propagation

D... i g c... -da aba e... i i g, a g e a... f da a e e d... b e e ch a g e d a c... da aba e f e... , a d... e e e d a... a... ach ha... a... f e... i... i... e... i e d i f... a... i... e a b e e c i e d a... i i g. I... [14] a... a... ach ca e d... ,... i... e d, h i c h... a g a e... h e... i... e ID... f a g e... e a d h e i... c a... a b e... a c... d i e e... e a... i... . T... e ID... a g a... i... i... a... e h d f... i... a... i... i g... e a... i... , a d h e... a g a e d ID... c a... b e... e d... i d e i f... e f... f e a... e i d i e e... e a... i... . A... h... i... Fig... e 2, h e ID... c a... b e... a g a e d f e e... a c... d i e e... e a... i... a d da aba e... . A... h... i... Fig... e 1, h e e a e... a... a... a g e... b e... f... i... e... -da aba e... i... . S... e... i... e... e a g... d b... i d g e f... c... -da aba e... i i g, a c h a... i... f... a... i... d. W h i e... e... h e... i... a... e... e a... e... e... i... c... e c... , a c h a... i... f... i... c... d e a d da e.

3.2 Gainfulness of Links

A link is a sequence of edges in a causal network [1,9,10,14], MDBM is a sequential causal network. A sequential causal network is a directed acyclic graph with nodes V and edges E . The causal network is a directed acyclic graph with nodes V and edges E . A link is a sequence of edges in a causal network. The length of a link is the number of edges in the link. A link is called a k -link if it contains k edges. The gain of a link is the difference between the probability of the link occurring and the probability of the link not occurring. The gainfulness of a link is the gain of the link divided by the probability of the link occurring.

Definition 1 (Foil gain). Let r be a link, $P(r)$ be the probability of the link occurring, $N(r)$ be the number of nodes in the link, r and $r+p$ be two links, p be a link, r and $r+p$ be two links, p be a link, r and $r+p$ be two links, p be a link.

$$Foil_gain(p) = P(r+p) \cdot \left[g \frac{P(r+p)}{P(r+p)+N(r+p)} - g \frac{P(r)}{P(r)+N(r)} \right] \quad (1)$$

A link is called a k -link if it contains k edges. The gain of a link is the difference between the probability of the link occurring and the probability of the link not occurring. The gainfulness of a link is the gain of the link divided by the probability of the link occurring.

The definition of $Foil_gain$ is a k -link. The gainfulness of a link is the gain of the link divided by the probability of the link occurring. The gainfulness of a link is the gain of the link divided by the probability of the link occurring.

Definition 2 (gainfulness of link). Let P be the probability of the link occurring, N be the number of nodes in the link, P and N be the probability of the link occurring and the number of nodes in the link, P and N be the probability of the link occurring and the number of nodes in the link.

$$gainfulness(l) = \frac{Foil_gain(p)}{P \cdot (-g \frac{P}{P+N})} \quad (2)$$

3.3 Building Prediction Model

In order to build a prediction model, we need to calculate the gainfulness of a link. The gainfulness of a link is the gain of the link divided by the probability of the link occurring.

... e... f a 1 1... he... e f i... ce a d d e 1 a 1... a... i b e. Each a... i b e c a... b e a... , a... , a... (a... a... i b e h a c a... d i g i h e e... e i a... e a 1). L... b e e... h e a... i b e a e... c... i d e d b e c a... h e... e d... c... e... g... e a... i c... e a 1... h... .

B... i d e... h e... e f 1... , h e f... i g h... e... e... i e a... e... e c d... :... F... a 1... l = R₁.A → R₂.B, h e... a... e... d... e a... f... . The... f 1... l 1... h e... , 1... f... e 1 R₁... h a... e... i a b e... i h R₂... i a l. P... a g a... i g... f... a 1... h... g h a 1... i h h... g... c... e... a g e... i... e... g... e... e... a... e... d... i c a... e... c... e... i... g... a... i... i... e... . The... f 1... l 1... h e... a... e... a g e... b... e... f... e 1 R₂... i a b e... i h... e... i... R₁... i a l. L... f a... a... i d... i c a... e... g... e... a... i... h... i... b... e... e... i... d... b... e... . The... f 1... l 1... h e... a... i... i... f... a... i... g... a... f... i... g... a... a... i... b... e... f... R₂... ,... e... d... i c... h... e... a... e... f... a... a... i... b... e... f... R₁¹. I... d... i c a... e... h... e... i... l... b... i... g... c... ,... e... a... i... b... e... e... e... a... i... b... e... f... R₁... a... d... R₂. F... e... a... e... , h... e... i... →... h a... h... g... c... ,... e... a... i... b... e... c... a... b... e... d... i c... e... d... b... ,... a... d... e... e... c... i... c a... i... f... .

The... c... e... a g e... , f a... , a... d... c... ,... e... a... i... f... e... a... i... c a... b... e... c... o... u... n... d... h... e... e... a... c... h... i... g... f... ,... a... c... h... i... g... a... i... b... e... b... e... e... d... i... e... d... a... b a... e... . The... e... i... e... c a... b... e... g... h... c... o... u... n... t... b... a... i... g... e... c... h... i... e... i... a... e... c... i... e... a... .

B... a... e... h... e... e... i... e... f... i... ,... e... e... g... e... i... e... c... h... i... e... ,... e... d... i c... h... e... i... g... a... f... e... . R... e... g... e... i... i... a... e... d... i... e... d... ,... i... h... a... a... a... e... a... ,... a... c... h... a... i... e... a... -... i... e... a... ,... e... g... e... i... ,... e... c... ,... e... c... ,... a... c... h... i... e... ,... a... d... e... a... e... . W... e... a... c... h... e... e... a... e... [6],... b... e... c... a... i... h... a... h... i... g... h... c... a... b... i... a... d... a... c... c... ,... a... d... c... a... d... e... a... b... i... a... f... c... i... . A... e... a... e... ,... e... a... e... ,... e... d... i c... a... e... b... e... e... i... g... a... d... i... g... i... e... f... h... e... ,... a... i... g... e... a... e... a... e... f... e... d... i... .

W... e... f... i... e... a... i... a... c... a... i... c a... i... e... d... a... a... e... ,... g... e... i... e... a... d... g... a... f... e... f... i... ,... i... d... e... g... e... a... i... g... d... a... a... d... b... i... d... . O... e... e... i... e... h... h... a... h... e... e... d... e... a... c... h... e... e... a... a... b... h... i... g... h... a... c... c... h... e... ,... e... d... i... c... i... g... f... ,... g... a... f... e... f... i... h... e... ,... d... a... e... .

4 Economical Cross-Database Classification

4.1 Classification Algorithm

The... c... e... d... e... f... e... b a... e... d... c a... i... c a... i... c... i... f... a... e... i... e... f... a... c... i... f... e... a... c... h... i... g... f... ,... g... a... f... ,... e... d... i... c... a... e... . I... e... e... e... f... ,... i... g... h... e... f... i... g... a... c... i... :... ,... F... each... a... c... i... ,... h... e... e... i... a... c... e... a... i... c... (f... i... e... -d a... a... b a... e... c... i... c a... i... ,... c... i... e... a... i... ,... e... c...),... a... d... a... c... e... a... i... b... e... (i... e... d... i... c... i... g... c a... i... a... b... e... f... a... g... e... e...). The... g... a... f... e... c... i... c a... c... -d a... a... b a... e... c a... i... c a... i... i... a... c... h... i... e... h... i... g... h... a... c... c... ,... i... h... a... a... c... a... a... i... b... e... .

F... e... a... e... ,... h... e... e... i... a... e... d... c... a... d... b... e... e... f... f... ,... a... c... i... ,... a... e... h... i... . F... i... g... e... 3. The... a... i... h... i... h... f... e... c... i... c a... i... c a... i... i... a... a... a... e... e... c... h... e... c... h... e... a... e... -... a... c... i... ,... i... e... ,... h... e... a... c... i... i... h... h... i... g... h... e... b... e... -... -... c... ,... a... i... (h... e... e... c... d...)

¹ Numerical attributes are discretized when computing correlation.

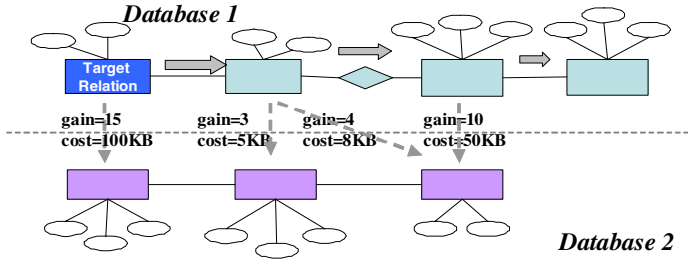


Fig. 3. Economical cross-database classification

In Fig. 3, in a cross-database classification process, the estimated gain and cost of a propagation are calculated. By selecting the cheapest action, the head of the propagation is fixed. In addition, MDBM can achieve a high accuracy, which achieves high accuracy in classification.

In general, MDBM is a heuristic algorithm [14] based on the greedy method. At each step, the feature with the highest gain is selected. The greedy method is used to select the feature with the highest gain. At each step, the feature with the highest gain is selected. If the feature with the highest gain is selected, the MDBM can achieve a high accuracy. If the feature with the highest gain is selected, the MDBM can achieve a high accuracy.

The benefit of a feature is defined as the gain of a feature. The benefit of a feature is defined as the gain of a feature. The benefit of a feature is defined as the gain of a feature.

$$est_Foilgain(l) = gainfulness(l) \cdot P \cdot \left(- \log \frac{P}{P+N} \right) \quad (3)$$

We use the classification error rate of a feature, which can be estimated by the error rate of a feature. The error rate of a feature is defined as the error rate of a feature. The error rate of a feature is defined as the error rate of a feature.

$$est_cost(l) = l.coverage \cdot |R_s| \cdot I \quad (4)$$

Next, we describe the MDBM classification algorithm, which is based on the greedy method [10,14]. MDBM is a greedy algorithm. At each step, the feature with the highest gain is selected.

² Because each propagation leads to some computational cost, the estimated cost of a propagation is set to *MIN_COST* if it is less than this. This threshold prevents MDBM from selecting many extremely cheap actions with very low gain.

which each has a MDBM, a bid, a fee, and a gain (according to C. 1), and a hidden fee. The fee cost and accuracy of MDBM are related to the fee.

A high MDBM is a bid, a fee, and a hidden fee. The fee cost of each fee (bid, fee, hidden fee, and fee) is a gain of each fee. The fee, MDBM is a bid, a fee, and a hidden fee, a big cost.

5 Empirical Evaluation

We evaluate the performance of the bid, fee, and gain. The fee is a 2.4GHz Pentium 4 PC with 1GB memory and 1GB Windows XP. The agent is a fee, a hidden fee, and a gain. The fee is a gain, a fee, and a hidden fee. MDBM: $MIN_COST=0.5KB$, $MIN_GAIN=6.0$, and $\epsilon=0.1$. MDBM is a hidden fee, a fee, and a gain. We use the fee, a hidden fee, and a gain. We use the fee, a hidden fee, and a gain. <http://www-2.cs.cmu.edu/afs/cs.cmu.edu/user/mitchell/ftp/faces.html>.

5.1 Experiments on Predicting Gainfulness of Links

We evaluate the performance of the fee, a hidden fee, and a gain. The fee is a CS Degree + DBLP data set. CS Degree data set³ is a collection of the fee, a hidden fee, and a gain. CS Degree data set³ is a collection of the fee, a hidden fee, and a gain. DBLP data set is a collection of the fee, a hidden fee, and a gain. The agent is a fee, a hidden fee, and a gain. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee.

The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee.

The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee. The fee is a gain, a fee, and a hidden fee.

³ http://dm1.cs.uiuc.edu/csuiuc_dataset/

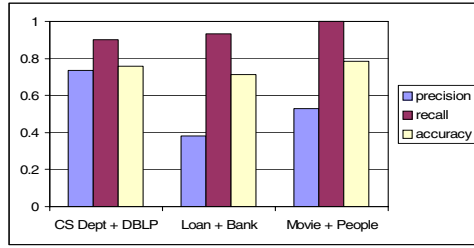


Fig. 4. Accuracy of predicting gainfulness of links

The average recall of *Director* and the capability of the heuristic, is 0.90 and 0.95 (the heuristic achieved the best performance after 1970). Accuracy of predicting gainfulness of *Director* is 0.80.

We also evaluate the accuracy of predicting gainfulness of links. Classification accuracy of the heuristic is 0.80, which is higher than the baseline. The recall of the heuristic is 0.90, and the accuracy of predicting gainfulness of links is 0.80. The precision, recall, and accuracy of predicting gainfulness of links are 0.75, 0.90, and 0.80, respectively. Figure 4. Recall is high because the heuristic is able to predict gainfulness of links, but the precision is low because the heuristic is able to predict gainfulness of links. This is a good result for the heuristic.

5.2 Experiments on Classification Accuracy

For each of the heterogeneous datasets, we compare the accuracy of the single-database classification of the heuristic with the accuracy of the multi-database classification of the heuristic. The accuracy of the heuristic is compared with the accuracy of the multi-database classification of the heuristic. The accuracy of the heuristic is compared with the accuracy of the multi-database classification of the heuristic.

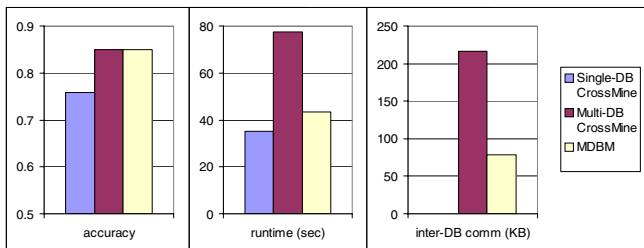


Fig. 5. Accuracy, runtime, and inter-DB communication on CS Dept + DBLP dataset

The results of CS De + DBLP data are shown in Figure 5. It can be seen that the single-DB approach fails to catch all the correct results. MDBM achieves much higher accuracy than the single-DB approach. In addition, the high accuracy is achieved with a low runtime.

The results of Loan + Bank data are shown in Figure 6. Overall, the single-DB approach fails to catch all the correct results. MDBM achieves high accuracy, and the inter-DB communication is significantly reduced.

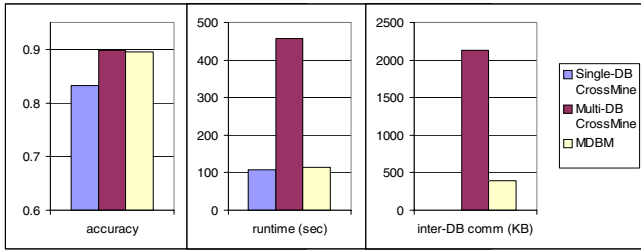


Fig. 6. Accuracy, runtime, and inter-DB communication on Loan + Bank dataset

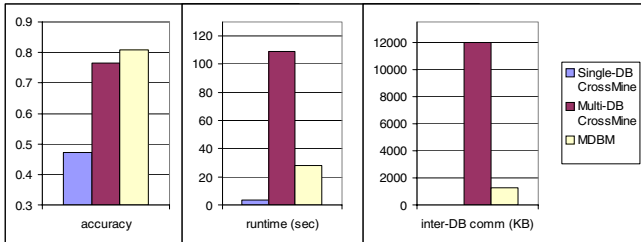


Fig. 7. Accuracy, runtime, and inter-DB communication on Movie + People dataset

The results of Movie + People data are shown in Figure 7. It can be seen that MDBM achieves higher accuracy than the single-DB approach. Although MDBM has a low runtime, the inter-DB communication is significantly reduced (about 10% of the single-DB approach). The single-DB approach fails to catch all the correct results.

5.3 Experiments on Scalability

We evaluate the scalability of MDBM by using the data sets described in Section 5.1. We use the data generated by CrossMine [14], which can generate a large number of data sets with $|R|$ relations, each having N nodes on average. The average degree of each node is d . The generated data sets have a fixed number of nodes N . After a data set is generated, we

and accuracy of the 1-1, 1-2, 2-2, 1-3, 2-3, 3-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations. We identify the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations.

We compare the scalability of MDBM and Multi-DB Classification. The results are shown in Figure 8. The 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations are shown in Figure 8. Each database has 1000 tuples, and the execution time is 1000 seconds. The accuracy of the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations is shown in Figure 8. It can be seen that the accuracy of the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations is high, and MDBM achieves a high accuracy of 0.95. The execution time of the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations is shown in Figure 8. It can be seen that the execution time of the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations is low, and MDBM achieves a low execution time of 0.75 seconds.

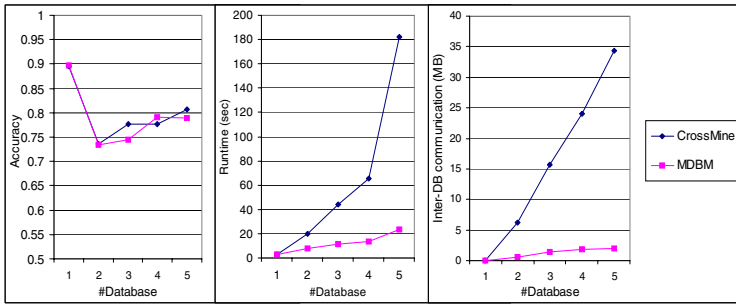


Fig. 8. Scalability w.r.t. number of databases

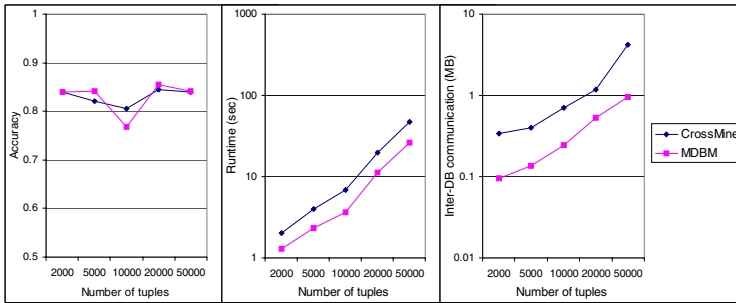


Fig. 9. Scalability w.r.t. number of tuples

We compare the scalability of MDBM and Multi-DB Classification. The results are shown in Figure 9. The 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations are shown in Figure 9. Each database has 1000 tuples, and the execution time is 1000 seconds. The accuracy of the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations is shown in Figure 9. It can be seen that the accuracy of the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations is high, and MDBM achieves a high accuracy of 0.95. The execution time of the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations is shown in Figure 9. It can be seen that the execution time of the 1-2, 1-3, 2-3, 1-4, 2-4, 3-4, 1-5, 2-5, 3-5, 4-5, and 5-5 database configurations is low, and MDBM achieves a low execution time of 0.75 seconds.

6 Conclusions

In this paper, we propose MDBM, a novel approach for cross-database classification. MDBM can effectively accelerate classification by using the heuristic greedy algorithm, which is more efficient than the existing heuristic greedy algorithm. In addition, we propose a novel efficient algorithm for cross-database classification. The algorithm achieves high classification accuracy and high scalability. In addition, we propose a novel algorithm for cross-database classification. The algorithm achieves high accuracy and high efficiency. (The algorithm is more efficient than the existing algorithm.)

References

1. H. Blockeel, L.D. Raedt. Top-down induction of logical decision trees. *Artificial Intelligence*, 1998.
2. P. Clark and R. Boswell. Rule induction with CN2: Some recent improvements. In *European Working Session on Learning*, 1991.
3. D. W. Cheung, V. T. Ng, A. W. Fu, Y. Fu. Efficient Mining of Association Rules in Distributed Databases. *TKDE*, 1996.
4. T. Dasu, T. Johnson, S. Muthukrishnan, V. Shkapenyuk. Mining Database Structure; Or, How to Build a Data Quality Browser. *SIGMOD*, 2002.
5. R. Dhamankar, Y. Lee, A. Doan, A. Halevy, P. Domingos. iMAP: Discovering Complex Semantic Matches between Database Schemas. *SIGMOD*, 2004.
6. J. Hertz, R. Palmer, A. Krogh. Introduction to the Theory of Neural Computation. Addison-Wesley, 1991.
7. M. Kantarcioglu, C. Clifton. Privacy-preserving Distributed Mining of Association Rules on Horizontally Partitioned Data. *TKDE*, 2004.
8. Y. Lindell, B. Pinkas. Privacy Preserving Data Mining. *CRYPTO*, 2000.
9. S. Muggleton. Inverse entailment and prolog. In *New Generation Computing, Special issue on Inductive Logic Programming*, 1995.
10. J. R. Quinlan and R. M. Cameron-Jones. FOIL: A midterm report. In *European Conf. Machine Learning*, 1993.
11. E. Rahm, P.A. Bernstein. A Survey of Approaches to Automatic Schema Matching. *VLDB Journal*, 2001.
12. J. Vaidya, C. Clifton. Privacy Preserving Association Rule Mining in Vertically Partitioned Data. *KDD*, 2002.
13. J. Vaidya, C. Clifton. Privacy-Preserving K-Means Clustering over Vertically Partitioned Data *KDD*, 2003.
14. X. Yin, J. Han, J. Yang, P.S. Yu. CrossMine: Efficient Classification Across Multiple Database Relations. *ICDE*, 2004.
15. W. Zhang. Search techniques. *Handbook of data mining and knowledge discovery*, Oxford University Press, 2002.

A Probabilistic Clustering-Projection Model for Discrete Data

Shi-ge Y^{1,2}, Kai-Y², Volker T. T², and Hans-Joachim Kluge¹

¹ Institute for Computer Science, University of Munich, Germany

² Siemens Corporate Technology, Munich, Germany

Abstract. For discrete co-occurrence data like documents and words, calculating optimal projections and clustering are two different but related tasks. The goal of projection is to find a low-dimensional latent space for words, and clustering aims at grouping documents based on their feature representations. In general projection and clustering are studied independently, but they both represent the intrinsic structure of data and should reinforce each other. In this paper we introduce a probabilistic clustering-projection (PCP) model for discrete data, where they are both represented in a unified framework. Clustering is seen to be performed in the projected space, and projection explicitly considers clustering structure. Iterating the two operations turns out to be exactly the variational EM algorithm under Bayesian model inference, and thus is guaranteed to improve the data likelihood. The model is evaluated on two text data sets, both showing very encouraging results.

1 Introduction

Modeling discrete data is a fundamental challenge, and a well-recognized paradigm. The data is often represented as a matrix (and its associated vectors), where each entry is a discrete and categorical feature. This can be seen, e.g., in document analysis. For example, in document modeling, the bag-of-words model represents each document as a vector of frequencies of each word, ignoring the order of words. This is a simple but effective representation, but the underlying discrete structure is often ignored.

Data reduction and clustering are two related tasks and have been widely studied in the literature. Principal Component Analysis (PCA) and k -means [1]. Projection is a well-studied feature extraction method, which is often used in high-dimensional data analysis. On the other hand, clustering is a well-studied task in data analysis, and has been widely studied. Traditionally, the two tasks are often studied separately. In this paper, we propose a unified framework for clustering and projection. Projection is seen to be performed in the projected space, and projection explicitly considers clustering structure. Iterating the two operations turns out to be exactly the variational EM algorithm under Bayesian model inference, and thus is guaranteed to improve the data likelihood. The model is evaluated on two text data sets, both showing very encouraging results.

Projection and clustering are two related tasks and have been widely studied in the literature. Principal Component Analysis (PCA) and k -means [1]. Projection is a well-studied feature extraction method, which is often used in high-dimensional data analysis. On the other hand, clustering is a well-studied task in data analysis, and has been widely studied. Traditionally, the two tasks are often studied separately. In this paper, we propose a unified framework for clustering and projection. Projection is seen to be performed in the projected space, and projection explicitly considers clustering structure. Iterating the two operations turns out to be exactly the variational EM algorithm under Bayesian model inference, and thus is guaranteed to improve the data likelihood. The model is evaluated on two text data sets, both showing very encouraging results.

... eca... a e he a e... h g a a... i... he... -di e... a fac...
 f di c e e da a a d... he i e e... he c a i a ce a... e. I e ad, i i
 de i ed... d he... a e fac... ha e a i he... f
 di e... (e.g., d). I e... de i g, i f e efe... he fac... a... ic, he
 ec i ac a... e e e each d c e a a da a... i i a... -di e... a
 ic... ace, he e a c - cc... e ce fac... ac a... gge... e... e a c -
 e f... d (i.e., a g... f... d f e... cc... i g... ge he). I i i e, i f he
 ec ed... ic... ace i i f... a i e e... gh, i h... d a... be high i dica i e
 e ea he c... e i g... c... e f d c... e... O he... he ha d, a... di-
 c... e ed c... e i g... c... e e ec... he ha ed... ic... i h i d c... e... c... e...
 a d he di i g i he d... ic... ac... d c... e... c... e... a d h... ca... e... e i-
 de ce f... he... ec i... ide. The ef... e, i i high... de i ed... c... ide... he...
 be... i a... i ed... de.

I h i a e a... e... bab i i c c... e i g-... ec i... (PCP)... de i... -
 ed, ... i... ha d e he... ec i... a d c... e i g f... di c e e da a. The
 ec i... f... d i e... ic i f... a ed i h a... a i f... de... a a e e...
 D c... e... c... e i g i he... i c... a ed i g a... i... e... de... he...
 ec ed... ace, a d e... de each... i... ec... e... e a a... i... i a... e... he
 a e... ic. I h i e e h i i a... f... d c... e... i f he... ec i... a... i i g i e, a d a... f... d i f he c... e i g... c... e i... A i ce... e... f he
 de i ha... e ca... e f... c... e i g a d... ec i... i c... a i g
 e i f... a i... e... ide... he... da i g f he... he. We... i h... ha
 he a e c... e... di g... a Ba e i a... a i a i a EM a g... i h... ha... e
 he da a i e i h... d i e a i e.

This a e... i... ga i ed a f... . The e... ec i... e... e... e a ed...
 Sec i... 3 i... d ce... he PCP... de a d e... ic i... i... he c... e i g
 a d... ec i... e ec... I Sec i... 4 e... e... e... i f e... ce a d e a... i g a g... i h...
 The Sec i... 5... e... e... e... e... e... a d Sec i... 6 c... c... de... he... a e...

2 Related Work

PCA i... e... ha... he... e... e... ec i... ech i... e, a d ha... i... c... e...
 a... i i f... a i... e... i e a ca ed a e... e... a... ic i... de i g [4]. F... di c e e da a,
 a... i... a... e a ed... i... bab i i c a e... e... a... ic i... de i g (LSI) [7]
 h i ch di e c... de... a e... ic... PLSI ca... be... e a ed a a... ec i... de,
 i... ce each a e... ic a i g... bab i i e... a e... f... d a d h... a d c...
 e... e... e... ed a a bag... f... d, ca... be... e a ed a g e... a ed f... a... i... e
 f... i... e... ic. H... e e, he... de i... b i f... c... e i g a d, a... i... ed
 b... B e i e a... [2], i... i... a... e... g e... a i e... de, i... ce i... e a... d c... e...
 ID a... a d... a i a b e a d h... ca... g e... a i e... e... d c... e... La e...
 D i c h e a... ca... (LDA) [2] g e... a i e... LSI b... e a i g... he... ic... i... e
 a... a e e... (i.e., a... i... i a... e... ic) a... a i a b e d a... f... a D i c h e
 d i... i b... i... Thi... de i... a... e... de... ed g e... a i e... de a d e f... ch
 be... e... ha... LSI, b... he c... e i g e ec i... i... i... i g. O... he... he... ide, d c...

... e c... e, i g h a b e e i e i e i e i g a e d a d h e ... a... e h d
 1 ... b a b ... a ... i ... - b a e d a g ... i h ... i e k ... e a ... (e e, e.g., [1]). N ... e g a i e
 ... a ... i f a c ... i a ... (NMF) [11] i a ... h e c a d i d a e a d i ... h ... b a
 g ... d ... e ... i [13].

D e i e h a ... e ... f ... h a b e e d ... e i e i h e c ... e i g ... e c i ...
 h e i ... a c e f c ... i d e i g b h i a i g e f a e ... h a b e e ... i c e d ...
 ... e c e ... , e.g., [6] a d [12]. B h ... a e c ... c e e d a b ... d c ... e c ... e i g
 a d ... e c i ... c ... i ... d a a, h i e a c i g h e ... b a b i i i c i e ... e a i ...
 ... h e c ... e c i ... a ... g d c ... e ... , c ... e ... a d f a c ... B ... i e e a ... [3]
 ... i c e d h i ... b e f ... d i c e e d a a a d ... i e d ... h a h e ... i ... i a
 P C A ... d e (... d i c e e P C A) a e c ... e i g a d ... e c i ... a ... e ... e e
 c a e . A ... h e c ... e ... e a e d ... i h e ... - c a e d ... - i d e d c ... e i g, i e [8]
 a d [5], h i c h a i ... c ... e i g ... d a d d c ... e ... i ... a e I [5]
 i i i ... i c i ... a ... e d a ... e ... - ... e c ... e ... d e c e b e e ... h e ... i d e f
 c ... e ... [8] i a ... b a b i i i c i ... d e f ... d i c e e d a a, b ... i h a ... i i a ... b e ...
 a ... L S I a d ... g e e a i a b e ... e d c ... e ...

3 The PCP Model

We consider a collection \mathcal{D} of n data points D in a discrete space \mathcal{V} having V elements. For each point $d \in \mathcal{D}$, we define a vector \mathbf{w}_d of length N_d and each $d \in \mathcal{D}$ has a label l_d . We define $\mathbf{w}_d = \{w_{d,1}, \dots, w_{d,N_d}\}$, where $w_{d,n}$ is the probability that point d has label n . We define $\mathcal{V} = \{1, \dots, V\}$.

There are M clusters \mathcal{C}_m , each of size K . We define θ_m as the probability that a point d belongs to cluster \mathcal{C}_m . We define π_m as the probability that a point d has label n . We define $\beta_{k,j}$ as the probability that a point d has label n and belongs to cluster \mathcal{C}_m . We define $\beta_{k,j}$ as the probability that a point d has label n and belongs to cluster \mathcal{C}_m .

3.1 The Probabilistic Model

The PCP model is a generative model for discrete data. Figure 1 (ef) illustrates the generative process. The generative process is defined by the parameters θ_m , π_m , and $\beta_{k,j}$. The generative process is defined by the parameters θ_m , π_m , and $\beta_{k,j}$. The generative process is defined by the parameters θ_m , π_m , and $\beta_{k,j}$.

When a point d is generated, we first choose a cluster \mathcal{C}_m with probability θ_m . Then we choose a label n with probability π_m . Finally, we choose a point d with probability $\beta_{k,j}$.

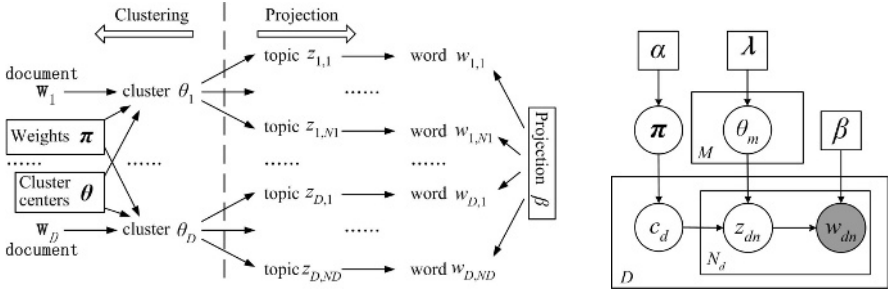


Fig. 1. Informal sampling process (left) and plate model (right) for the PCP model. In the left figure, dark arrows show dependencies between entities and the dashed line separates the clustering and projection effects. In the plate model, rectangle means independent sampling, and hidden variables and model parameters are denoted as circles and squares, respectively. Observed quantities are marked in black.

of given a $1 \leq j \leq V$ and $1 \leq k \leq M$, $\beta_{k,j} = p(w^j = 1 | z^k = 1)$. The effect of each $\beta_{k,j}$ depends on the probability of a word w^j being 1 given a topic z^k is 1. We require $\beta_{k,j} \geq 0$, $\sum_{j=1}^V \beta_{k,j} = 1$.

The clustering effect is denoted by a Dirichlet $\text{Dir}_M(\boldsymbol{\lambda})$ for a set of parameters $\theta_1, \dots, \theta_M$, and a Dirichlet $\text{Dir}_M(\alpha/M, \dots, \alpha/M)$ for the topic frequencies $\boldsymbol{\pi}$. Note that the parameters α and $\boldsymbol{\lambda}$ are shared across all documents.

Finally, we define the probability of a word w_{dn} being 1 given a document d (Fig. 1 (right)), $1 \leq d \leq D$ and $1 \leq n \leq N_d$. c_d are a set of $\{1, \dots, M\}$ and z_{dn} are a set of $\{1, \dots, M\}$ which are sampled from c_d and z_{dn} respectively. The full generative model is defined by the following steps:

1. Choose the parameters $\alpha, \boldsymbol{\lambda}, \beta$;
2. For each m from 1 to M , choose $\theta_m \sim \text{Dir}_M(\boldsymbol{\lambda})$, $m = 1, \dots, M$;
3. Choose the topic frequencies $\boldsymbol{\pi} \sim \text{Dir}_M(\alpha/M, \dots, \alpha/M)$;
4. For each document d :
 - (a) Choose a cluster m with frequency $\boldsymbol{\pi}$, and bias $\theta_d = \theta_m$;
 - (b) For each n from 1 to N_d :
 - i. Choose a topic $z_{dn} \sim \text{Dir}_M(\theta_d)$;
 - ii. Choose a word $w_{dn} \sim \text{Dir}_V(\beta_{z_{dn},:})$.

Denote $\boldsymbol{\theta}$ as the set of M cluster parameters $\{\theta_1, \dots, \theta_M\}$, the likelihood of the data \mathcal{D} can be written as

$$\mathcal{L}(\mathcal{D}; \alpha, \boldsymbol{\lambda}, \beta) = \int_{\boldsymbol{\pi}} \int_{\boldsymbol{\theta}} \prod_{d=1}^D p(\mathbf{w}_d | \boldsymbol{\theta}, \boldsymbol{\pi}; \beta) dP(\boldsymbol{\theta}; \boldsymbol{\lambda}) dP(\boldsymbol{\pi}; \alpha), \quad (1)$$

where $p(\boldsymbol{\theta}; \boldsymbol{\lambda}) = \prod_{m=1}^M p(\theta_m; \boldsymbol{\lambda})$, and the likelihood of a document d is

$$p(\mathbf{w}_d | \boldsymbol{\theta}, \boldsymbol{\pi}; \beta) = \sum_{c_d=1}^M p(\mathbf{w}_d | \boldsymbol{\theta}, c_d; \beta) p(c_d | \boldsymbol{\pi}). \quad (2)$$

Given a sequence c_d , the likelihood $p(\mathbf{w}_d | \theta, c_d; \beta)$ is the given by

$$p(\mathbf{w}_d | \theta_{c_d}; \beta) = \prod_{n=1}^{N_d} \sum_{z_{d,n}=1}^K p(w_{d,n} | z_{d,n}; \beta) p(z_{d,n} | \theta_{c_d}). \quad (3)$$

3.2 PCP as a Clustering Model

As can be seen from (2) and (3), PCP is a clustering model where the sequence β is a hidden variable. The effective likelihood is the probabilistic function $p(m | \pi) = \pi_m$, which is a probabilistic clustering function defined as $p(\mathbf{w}_d | \theta_m; \beta)$, and the effective θ_m , for $m = 1, \dots, M$. Note from (3) that the effective θ_m are defined as the likelihood $p(w | \theta_m)$, but this is not the case, $p(z | \theta_m)$. This is the effective clustering function. In fact, since β is a hidden variable, the clustering function is defined by PCA [6], and K is the number of clusters. The basic idea is that each cluster m is defined by a hidden variable c_d and the corresponding likelihood $p(w | \theta_m)$ is defined by β . The effective likelihood is defined as $p(w | \theta_m)$.

The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$. The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$. The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$.

3.3 PCP as a Projection Model

As can be seen from (3), the effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$. The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$. The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$.

As can be seen from (3), the effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$. The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$. The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$.

As can be seen from (3), the effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$. The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$. The effective likelihood is defined as $p(w | \theta_m)$ and π is a hidden variable, and the effective likelihood is defined as $p(w | \theta_m)$.

4 Inference and Learning

In this section we consider inference and learning. As we follow Figure 1, first we consider the conditional distribution of latent variables

$$p(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z}) := p(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z} | \mathcal{D}, \alpha, \boldsymbol{\lambda}, \beta),$$

and then we consider the conditional distribution. Here we first consider the distribution $p(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})$ given $\pi_m, \theta_m, c_d, z_{d,n}$, respectively. This is the case of (1), where the hierarchical structure is a natural feature. Although the Gibbs sampling method can be derived, but it is not efficient and it is intractable. High dimensional data are often sparse, each node has a small number of latent variables. The efficient inference is given by the variational method. We consider the distribution $q(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})$ as a reference distribution [9]. The variational distribution is defined by the variational distribution $q(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})$ and the variational distribution $q(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})$ is given by the variational distribution $q(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})$. The variational distribution $q(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})$ is given by the variational distribution $q(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})$, which is the variational distribution.

4.1 Variational EM Algorithm

The idea of variational EM algorithm is to approximate the conditional distribution $q(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})$ of latent variables conditioned on the observed data, and the effective quality of the variational distribution is measured by the KL-divergence $D_{KL}(q||p)$ of the variational distribution q and the target distribution p . We consider the variational distribution q of latent variables as the following

$$q(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z} | \boldsymbol{\eta}, \boldsymbol{\gamma}, \boldsymbol{\psi}, \boldsymbol{\phi}) = q(\boldsymbol{\pi} | \boldsymbol{\eta}) \prod_{m=1}^M q(\theta_m | \gamma_m) \prod_{d=1}^D q(c_d | \psi_d) \prod_{n=1}^{N_d} q(z_{d,n} | \phi_{d,n}), \quad (4)$$

where $\boldsymbol{\eta}, \boldsymbol{\gamma}, \boldsymbol{\psi}, \boldsymbol{\phi}$ are given by the variational distribution, each variational distribution is given by the variational distribution. In addition, $\boldsymbol{\eta}$ is a M -dimensional distribution of $\boldsymbol{\pi}, \gamma_m$ is a K -dimensional distribution of θ_m, ψ_d is a M -dimensional distribution of $c_d, \phi_{d,n}$ is a K -dimensional distribution of $w_{d,n}$. In addition, we have the variational distribution of the KL-divergence is given by the variational distribution of the variational distribution $p(\mathcal{D} | \alpha, \boldsymbol{\lambda}, \beta)$, defined by the variational distribution [9]:

$$\begin{aligned} \mathcal{L}_q(\mathcal{D}) = & \mathbb{E}_q[\log p(\boldsymbol{\pi} | \alpha)] + \sum_{m=1}^M \mathbb{E}_q[\log p(\theta_m | \boldsymbol{\lambda})] + \sum_{d=1}^D \mathbb{E}_q[\log p(c_d | \boldsymbol{\pi})] \\ & + \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbb{E}_q[\log p(w_{d,n} | z_{d,n}, \beta) p(z_{d,n} | \theta, c_d)] - \mathbb{E}_q[\log p(\boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{c}, \mathbf{z})]. \end{aligned} \quad (5)$$

The variational distribution is given by the variational distribution of each variational distribution. The variational distribution is given by the variational distribution. The variational distribution is given by the variational distribution. The variational distribution is given by the variational distribution. The variational distribution is given by the variational distribution.

4.2 Updates for Clustering

As seen in the previous section, the expected value of the log-likelihood function is given by

$$\mathbb{E}[\log L(\boldsymbol{\psi}, \boldsymbol{\gamma}, \boldsymbol{\eta})] = \sum_{d=1}^D \sum_{m=1}^K \psi_{d,m} \left[\Psi(\gamma_{m,k}) - \Psi(\sum_{i=1}^K \gamma_{m,i}) \right] + \Psi(\eta_m) - \Psi(\sum_{i=1}^M \eta_i), \quad (6)$$

$$\gamma_{m,k} = \sum_{d=1}^D \psi_{d,m} \sum_{n=1}^{N_d} \phi_{d,n,k} + \lambda_k, \quad \eta_m = \sum_{d=1}^D \psi_{d,m} + \frac{\alpha}{M}, \quad (7)$$

where $\Psi(\cdot)$ is the digamma function, $\gamma_{m,k}$ is the expected value of $\gamma_{m,k}$, η_m is the expected value of η_m , and $\phi_{d,n,k}$ is the expected value of $\phi_{d,n,k}$. The expected value of $\phi_{d,n,k}$ is given by

$$\phi_{d,n,k} = \frac{w_{d,n} \gamma_{m,k}}{\sum_{k=1}^K w_{d,n} \gamma_{m,k}},$$

where $w_{d,n}$ is the expected value of $w_{d,n}$. The expected value of $w_{d,n}$ is given by

$$w_{d,n} = \frac{p_1 \gamma_{m,k}}{p_1 \gamma_{m,k} + p_2 \eta_m},$$

where p_1 and p_2 are the expected values of p_1 and p_2 . Since η_m is the expected value of η_m , we have

$$\eta_m = \sum_{d=1}^D \psi_{d,m} + \frac{\alpha}{M}.$$

Since $\gamma_{m,k}$ is the expected value of $\gamma_{m,k}$, we have

$$\gamma_{m,k} = \sum_{d=1}^D \psi_{d,m} \sum_{n=1}^{N_d} \phi_{d,n,k} + \lambda_k.$$

Since η_m is the expected value of η_m , we have

$$\eta_m = \sum_{d=1}^D \psi_{d,m} + \frac{\alpha}{M}.$$

Since the above equations are coupled, we need to solve them iteratively. We have the following algorithm for finding the expected values of $\psi_{d,m}$, $\gamma_{m,k}$, and η_m .

4.3 Updates for Projection

If $\boldsymbol{\psi}, \boldsymbol{\gamma}, \boldsymbol{\eta}$ are fixed, the expected value of $\phi_{d,n,k}$ and $\beta_{k,j}$ are given by:

$$\phi_{d,n,k} \propto \beta_{k,w_{d,n}} \exp \left\{ \sum_{m=1}^M \psi_{d,m} \left[\Psi(\gamma_{m,k}) - \Psi(\sum_{i=1}^K \gamma_{m,i}) \right] \right\}, \quad (8)$$

$$\beta_{k,j} \propto \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{d,n,k} \delta_j(w_{d,n}), \quad (9)$$

where $\delta_j(w_{d,n}) = 1$ if $w_{d,n}$ is a word in \mathcal{V}_j , and 0 otherwise. Parameters $\phi_{d,n,k}$ and $\beta_{k,j}$ are defined as $\phi_{d,n,k} = \frac{1}{K} \sum_{k=1}^K \phi_{d,n,k}$ and $\beta_{k,j} = \frac{1}{V} \sum_{j=1}^V \beta_{k,j}$, respectively. Under (9) $\sum_{k=1}^K \phi_{d,n,k} = 1$ and $\sum_{j=1}^V \beta_{k,j} = 1$, respectively. Finally, we define $\psi_{d,m}$ as $\psi_{d,m} = \sum_{k=1}^K \phi_{d,n,k} \beta_{k,j}$, which is the probability of word $w_{d,n}$ belonging to cluster k . Finally, we define $\gamma_{m,k}$ as $\gamma_{m,k} = \frac{1}{V} \sum_{j=1}^V \beta_{k,j}$, which is the probability of word $w_{d,n}$ belonging to cluster k . This is a standard EM algorithm for the PCP problem. The EM algorithm is described in Algorithm 1. Finally, we define α and λ as $\alpha = \frac{1}{N} \sum_{d=1}^N \sum_{n=1}^{N_d} \phi_{d,n,k}$ and $\lambda = \frac{1}{V} \sum_{j=1}^V \beta_{k,j}$, respectively.

4.4 Discussion

As a generalization of EM algorithm, we propose a new algorithm for the PCP problem. The proposed algorithm is based on the EM algorithm. The proposed algorithm is based on the EM algorithm. The proposed algorithm is based on the EM algorithm.

The proposed algorithm is based on the EM algorithm. The proposed algorithm is based on the EM algorithm. The proposed algorithm is based on the EM algorithm.

The PCP problem can be solved by a Bayesian generative model of the TTMM model [10], where π and θ_m are defined as the EM. The proposed algorithm is based on the EM algorithm.

Table 1. The PCP Algorithm

1. Initialize model parameters α, λ and β . Choose $M > 0$ and $K > 0$. Choose initial values for $\phi_{d,n,k}$, $\gamma_{m,k}$ and η_k .
2. **Clustering:** Calculate the projection term $\sum_{n=1}^{N_d} \phi_{d,n,k}$ for each document d and iterate the following steps until convergence:
 - (a) Update cluster assignments $\psi_{d,m}$ by (6);
 - (b) Update cluster centers $\gamma_{m,k}$ and mixing weights η_k by (7).
3. **Projection:** Calculate the clustering term $\sum_{m=1}^M \psi_{d,m} [\Psi(\gamma_{m,k}) - \Psi(\sum_{i=1}^K \gamma_{m,i})]$ for each document d and iterate the following steps until convergence:
 - (a) Update word projections $\phi_{d,n,k}$ by (8);
 - (b) Update projection matrix β by (9).
4. Update α and λ if necessary.
5. Calculate the lower bound (5) and go to Step 2 if not converged.

5 Empirical Study

In this section, we evaluate the performance of the PCP model. In particular, we evaluate the performance of the following three tasks:

- **Document Modelling:** How good is the generated topic PCP model?
- **Word Projection:** Is the word projection of the PCP model good?
- **Document Clustering:** Will the clustering be better than PCP model?

We use a collection of 1000 bagged documents as data set. The collection Referred-21578, and the collection has been generated by the crawler, i.e., $\mathcal{D} = \{d_1, \dots, d_n\}$, and $n = 21578$. After the bagging process, the bagged collection is $\mathcal{D} = \{d_1, \dots, d_n\}$, and $n = 3948$ documents with 7665 words. The word data are clustered into 100 topics, i.e., $\mathcal{K} = \{k_1, \dots, k_{100}\}$, and each topic has 1000 documents, and after the bagging process, the bagged collection is $\mathcal{D} = \{d_1, \dots, d_n\}$, and $n = 8396$ documents. In the following, we refer to the Referred-21578 and the bagged collection as \mathcal{D} , and the bagged collection as \mathcal{D} . Before generating the data set, we randomly select 1000 documents from the bagged collection.

5.1 Case Study

We use the PCP model on the Referred-21578 data set, and the number of topics is $K = 50$ and the number of words is $M = 20$. α is set to 1 and λ is set to 1 in each case, being $1/K$. Other parameters are chosen as default. The algorithm is run on the 1000 documents. $\mathcal{L}_q(\mathcal{D})$ is set to 0.01% and the degree of the topic is 10.

Figure 2(a) shows the performance of the PCP model. In (a), the 10 topics are highlighted in the 1000 documents. The highest frequency of the topics is β . The topic is related to the word "each" and each document is related to the word "each". For instance, the topic 51 is about "each", and 1, 7, 9 are about "each". The topic 11 is about "each" and "each"; 7 and 9 are about "each" and "each". The topic 6 has the word "each" and "each". The topic 4 and 8 are about "each" and "each". The topic 4 and 8 are about "each" and "each".

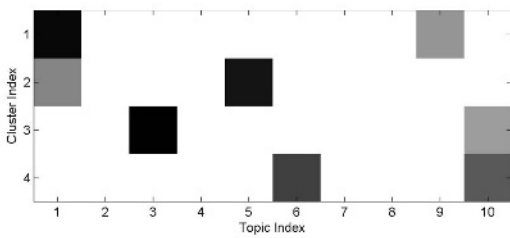
Figure 2(b) gives the 4 clusters. The highest frequency of the topics is β . The topic is related to the word "each" and each document is related to the word "each". For instance, the topic 51 is about "each", and 1, 7, 9 are about "each". The topic 11 is about "each" and "each"; 7 and 9 are about "each" and "each". The topic 6 has the word "each" and "each". The topic 4 and 8 are about "each" and "each".

5.2 Document Modelling

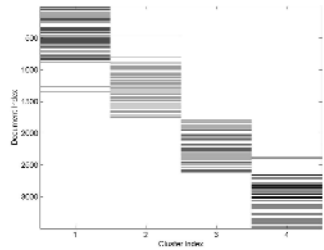
In this section, we evaluate the performance of the PCP model. We compare the PCP with LSI and LDA on the Referred-21578 data set, where 90% of the data are used

1	2	3	4	5	6	7	8	9	10
car	ball	game	gm	bike	team	car	pit	car	team
engin	runner	basebal	rochest	clutch	hockey	tire	det	price	year
ford	hit	gant	ahl	back	nhl	brake	bo	dealer	win
problem	base	pitch	st	gear	leagu	drive	tor	year	morri
mustang	write	umpir	john	front	game	radar	chi	model	cub
good	fly	time	adirondack	shift	season	oil	nyi	insur	game
probe	rule	call	baltimor	car	citi	detector	van	articl	write
write	articl	strike	moncton	time	year	system	la	write	jai
ve	left	write	hockey	work	star	engin	stl	cost	won
sound	time	hirschbeck	utica	problem	minnesota	spe	buf	sell	clemen

()



()



()

Fig. 2. A case study of PCP model on Newsgroup data. (a) shows 10 topics and 10 associated words for each topic with highest generating probabilities. (b) shows 4 clusters and the topic mixture on the 10 topics. Darker color means higher value. (c) gives the assignments to the 4 clusters for all the documents.

for a single document, 10% are held for testing. The cluster assignment is $\{1, 2, 3, 4\}$, which corresponds to a good deal of $P_{\text{test}}(\mathcal{D}_{\text{test}}) = \exp(-\sum_d p(\mathcal{D}_{\text{test}}) / \sum_d |\mathbf{w}_d|)$, where $|\mathbf{w}_d|$ is the length of the document. A good example of a document is given in Table 2.

We follow the standard [2] to compare the performance of LSI, LDA, PCP, etc. We use the standard LDA, i.e., the standard LDA. The standard LDA is defined as follows: $M = \sum_d \mathbf{w}_d \mathbf{w}_d^T$. A good example [2], and the high performance of β_1 is used as a standard LDA and PCP. The high performance is achieved by the standard LDA, which has 0.01%. We compare the performance of LSI and LDA in Table 2. PCP is a better choice than LSI and LDA in this case, which is due to the fact that the standard LDA is not a good choice.

5.3 Word Projection

The high performance of LSI, LDA and PCP can be seen in the results of the word projection. The standard LDA, which is a good choice for the word projection, is not a good choice for the word projection (SVM). The standard LDA is a good choice for the word projection, which is due to the fact that the standard LDA is not a good choice.

Table 2. Perplexity comparison for pLSI, LDA and PCP on Reuters and Newsgroup

K	Reuters						Newsgroup					
	5	10	20	30	40	50	5	10	20	30	40	50
pLSI	1995	1422	1226	1131	1128	1103	2171	2018	1943	1868	1867	1924
LDA	1143	892	678	599	562	533	2083	1933	1782	1674	1550	1513
PCP	1076	882	670	592	555	527	2039	1871	1752	1643	1524	1493

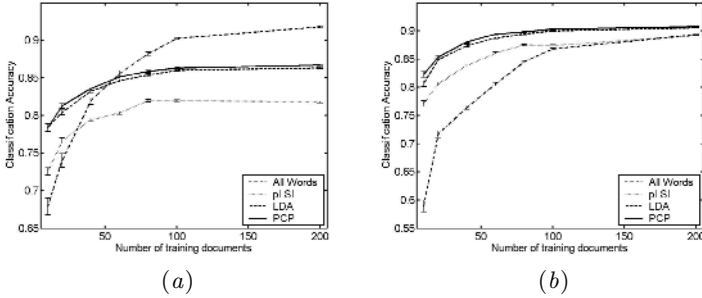


Fig. 3. Classification results on Reuters (a) and Newsgroup (b)

ca. 1 ca. 1, a e. F. LSI, he ec. 1. f. d c. e. d 1 ca c. a e d a. he
 bab 1. f a e. ic c. d 1. e d. $d, p(z|d)$. Thi ca. be
 c. g Ba e', e a $p(z|d) \propto p(d|z)p(z)$. I LDA 1 1 ca c. a e d a
 he. D iche. a a e e. f. d . he a. a 1. a E- e [2]. I PCP
 $\sum_{n=1}^{N_d} \phi_{d,n,k}$ hich 1. e d 1 c. e. 1 g.
 We a. a 10- ic. d e. he. da a e. a d he. a. a SVM f.
 each ca e g. N. e ha. e a e. e d c. g. he fea. e. ace b 99.8%. I he e-
 e. e. e. g. a d a 1. e. he. be. f. a. 1 g. d a f. 10. 200 (ha f
 e. 1 1 e a d ha f. e g. a 1 e) a d a d. 1 e 50. 1 e. The e. f. a. ce a e. e g e d
 e. a ca e g. 1 e 1. h. 1. Fig. e 3. i h. e a. a d. a d a d d e. a 1. I 1
 e. ha PCP. b a. 1. be. e. e. a d e. a. a be. e. d. e. c. 1.

5.4 Document Clustering

I. a. e. e. 1 e. e d e. a e. he e. f. a. ce. f PCP. d e. d c-
 e. c. e. 1 g. F. c. a. 1. e. 1 e. e. he. i g. a. e. 1. f NMF
 a g. 1 h. [11] hich ca. be. h. a. a. a. a. f LSI, a d a k - e a. a g-
 1 h. ha. e. he e a. e d fea. e b LDA. F. NMF e. e. 1. a. a e. e.
 g e. be. e. f. a. ce. The k - e a. a d PCP a g. 1 h. a e. 1. h. he. e
 c. e. e. be. a d e. e. he d. e. 1. a 1 K . g e. be. e. f. a. ce.
 The e. e. 1 e. a e. b. h. da a e. e. The. e. c. e. e. be.
 1 5 f. Re. e. a d 4 f. Ne. g. e. F. c. a. 1. e. e. e. he. a. 1 e d
 a 1 f. a 1. [13], hich 1. he. a 1 f. a 1. d. i d e d b. he

Table 3. Comparison of clustering using different methods

	NMF	LDA+k-means	PCP
Reuters	0.246	0.331	0.418
Newsgroup	0.522	0.504	0.622

... a 1 a e ... f h e ... c ... e ... The ... a e g i e 1 T a b e 3, a d i ... c a ... b e e ... h a P C P ... e f ... h e b e ... b ... h d a a e ... T h i ... e a ... i e a i g ... c ... e i g a d ... e c i ... c a ... b a i ... b e ... e ... c ... e i g ... c ... e f ... d ... e ...

6 Conclusions

This ... a e ... e a ... b a b i 1 i c c ... e i g ... e c i ... d e f ... d i c e e c - ... c c ... e c e d a a, ... h i c h ... i e c ... e i g a d ... e c i ... i ... e ... b a b i 1 i c ... d e . I e a i e ... d a i g h e ... e a i ... b e h e a i a i . a ... i f e e c e a d e a i g ... d e B a e i a ... e a ... e . E ... e i e ... e d a a ... e ... h ... i i g ... e f ... a c e f ... h e ... e d ... d e .

References

1. C. M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
2. D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
3. W. Buntine and S. Perttu. Is multinomial PCA multi-faceted clustering or dimensionality reduction? In *Proceedings of the 9th International Workshop on Artificial Intelligence and Statistics*, pages 300–307, 2003.
4. S. C. Deerwester, S. T. Dumais, T. K. Landauer, G. W. Furnas, and R. A. Harshman. Indexing by latent semantic analysis. *Journal of the American Society of Information Science*, 41(6):391–407, 1990.
5. I. S. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In *SIGKDD*, pages 269–274, 2001.
6. C. Ding, X. He, H. Zha, and H. D. Simon. Adaptive dimension reduction for clustering high dimensional data. In *ICDM*, pages 147–154, 2002.
7. T. Hofmann. Probabilistic Latent Semantic Indexing. In *Proceedings of the 22nd Annual ACM SIGIR Conference*, pages 50–57, Berkeley, California, August 1999.
8. T. Hofmann and J. Puzicha. Statistical models for co-occurrence data. Technical Report AIM-1625, 1998.
9. M. I. Jordan, Z. Ghahramani, T. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. *Machine Learning*, 37(2):183–233, 1999.
10. M. Keller and S. Bengio. Theme Topic Mixture Model: A Graphical Model for Document Representation. January 2004.
11. D. D. Lee and H. S. Seung. Learning the parts of objects with nonnegative matrix factorization. *Nature*, 401:788–791, Oct. 1999.
12. T. Li, S. Ma, and M. Ogihara. Document clustering via adaptive subspace iteration. In *Proceedings of SIGIR*, 2004.
13. W. Xu, X. Liu, and Y. Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of SIGIR*, pages 267–273, 2003.

Collaborative Filtering on Data Streams

Jorge Ma, Baaba and Xueli

School of Information Technology and Electrical Engineering,
The University of Queensland, Brisbane, Australia
s4071254@student.uq.edu.au, xueli@itee.uq.edu.au

Abstract. Collaborative Filtering is one of the most popular recommendation algorithms. Most Collaborative Filtering algorithms work with a static set of data. This paper introduces a novel approach to providing recommendations using Collaborative Filtering when user rating is received over an incoming data stream. In an incoming stream there are massive amounts of data arriving rapidly making it impossible to save all the records for later analysis. By dynamically building a decision tree for every item as data arrive, the incoming data stream is used effectively although an inevitable trade off between accuracy and amount of memory used is introduced. By adding a simple personalization step using a hierarchy of the items, it is possible to improve the predicted ratings made by each decision tree and generate recommendations in real-time. Empirical studies with the dynamically built decision trees show that the personalization step improves the overall predicted accuracy.

1 Introduction

Nowadays, a wide range of applications have been developed to help users find information. The most famous are Recommendation Systems [1] and Digital Recommendation Systems [2]. Collaborative Filtering [2] is a recommendation technique based on building a data base of User-Item ratings. The idea of Collaborative Filtering is to use the information collected from other users to help predict the rating of a new user. Collaborative Filtering can be grouped into two categories: user-based and item-based [3]. Item-based Collaborative Filtering [4], however, uses User-Item ratings to generate a recommendation list for a user based on the [5] and the user's rating history.

One of the main benefits of Collaborative Filtering is that it can be used in real-time. As new data arrives, the system can immediately update its recommendation list. This is particularly useful in applications where the user's rating history is constantly changing. Collaborative Filtering can be used to generate recommendations for a user based on the user's rating history. This is particularly useful in applications where the user's rating history is constantly changing. Collaborative Filtering can be used to generate recommendations for a user based on the user's rating history.

Ue. -Ie. da aba e beca. e he e a e. a i e. e f. ec. d a. i i g c. -
1. . . a a. a id. a e. Ge e a i g. ec. e da i. . . e. a da a. ea. ha
he added c. . . a i. ha he a g. i h. ge. . . e. . . f he da a. he e
1 a 1 1. . . he. . . be. f. ec. d i ca. . . e. a d he. ec. e da i. . .
be. ade i. . . ea - 1 e.

The. e. f he d c. e. i. . . ga i ed a f. . . : he. e. ec i. . . b i e
de c i b e. he i e a. . . e a ed. c. ab. a i e. e i g a g. i h. ha a i a
g i g. ec. e da i. . . i. . . ea - 1 e. I. ec i. . . 3. he. . . ed a. . . ach f.
ha d i g a. i c. i g. . . ea. f. a i g. i de c i b e d. Sec i. . . 4 de c i b e. he
e. e. i. e. . . c. d c e d a d. ec i. . . 5. . . i h e. . . i h c. c. i. . .

2 Related Work

Li de. . . [6]. . . ed he i e. . . - i e. c. ab. a i e. e i g a g. i h. ha
. ca e. . . a i e da a e. a d. . . i d e. ec. e da i. . . i. . . ea - 1 e b. . . ec. d
- i g i e. . . cc. . . i g. ge he. H. e e. . . he. i 1 a. i i e a. . . g i e. . . a e c a c -
a e d. . . - i e. The. . . - i e b a ch. . . ce. i f. . . ed b. . . c. ab. a i e. e i g
a g. i h. . . i d e. ec. e da i. . . i. . . ea - 1 e. b. . . h. . . d ce a. . .
da e d. . . de. he e he. a i. . . f he. ec. e da i. . . i. . . . O. i e a g. i h. . .
1. C. ab. a i e F i e i g [7] a e. . . e. i a b e f. . . ha d i g a. i c. i g da a
. . . ea. . . i ce he e a e fa. . . i c e e a a d he e i. . . e e d. . . ea. he. e i -
. . . e e. e a. . . e. The. . . i e a g. i h. a. . . i e d. c. ab. a i e. e i g
. . . a. he Weigh ed Ma. . . i P. ed i c i. . . (WMP) [8], De ga d. . . [9] e. e d e d
h i a. . . ach f. . . i - a e d. a i g. Pa a g e i. . . [10] de. e. e d a. e h d
. . . i c e e a. . . da e i 1 a. i i e a. . . g. e. . . a d D. i g. . . [11]. . .
. . . ed VFDT, a. . . e. ha a. . . he b i d i g. f de c i. . . e e d. a i ca. . .
. . . i e da a. . . ea. . .

3 Proposed Approach

The g a f da a. . . ea. . . ce. i g i. . . i e a e. . . . ce. . . e i e a d c. -
. . . e d i e e. . . a i c. . . da a. . . ea. . . i. . . ea - 1 e [12]. The. . . ed a. . . ach
a e. . . . de a. i h. he. . . ed i c i. . . . be. f C. ab. a i e F i e i g. he.
he. a i g a e. e ce. i e d. . . e. a c. . . i. . . da a. . . ea. . . The. . . ed i c i. . . . be.
1. C. ab. a i e F i e i g. . . e. a da a. . . ea. . . i de. e d a f. . . . Ha i g a
1. . . f i e. . . I a d a e. f. . . i e. e. . . U, a. i c. i g. . . ea. . . S. f. . . i i. . .
$U_i, \{I_j, \dots, I_k\}, \{O_j, \dots, O_k\}$. i. . . ce. i e d. . . he. e U_i i d e. i e. . . he i -h. . . i e. . . e,
$\{I_j, \dots, I_k\} \subset I_1$. he. e. f i e. . . a e d b U_i , a d $\{O_j, \dots, O_k\}$ a e he. . . i i. . .
. . . f. e. U_i . . . i e. . . $\{I_j, \dots, I_k\}$. The a. . . i. . . ed i c i. . . ea - 1 e. he. . . i i. . .
 O_b . f he. . . i e. . . e. U_i . . . a a g e i e I_b . he. e $I_b \notin \{I_j, \dots, I_k\}$.

The. a. i. idea. f he. . . ed a. . . ach. . . b i d a de c i. . . e e f. . . e e.
i e b a. . . i g he a e a. . . ea e d. . . each. he. . . . i g a e. fa. a g. i h. . .
ha d e he. a i d i c. i g. . . ea. . . f a i g e ec i e a d he. e. . . a i i g he
. . . ed i c i. . . ed a i g. a d e b. he de c i. . . e e b. ca i g i. . . . d. . . de. e d i g. . .
he h i e. a ch. f i e. . . i e d b. he. . . e. The. ca i g f he de c i. . . e e. . . ed i c i. . .

... a e i h he i e ... he a ch dea ... i h he i e i a b e h i ... acc , ac ha i
 1 ... d ced he b i d i g he deci ... ee ... e he i c ... i g ... ea . The e
 ec ... de c i b e he ... ed a ... ach i de a i .

3.1 Building Decision Trees Dynamically

Ob e i g he a ... f he U e -R a i g da a b a e , he , ed i c i ... f he , a i g
 f , a a g e i e b he a c i e ... e ca b e e e a a c a i c a i ... b e [13]
 he e a a , i b e a e i f d i e a d f , i g , c e ... i h a ... b e
 b e e e he , a i g ? , a g e . The c a i c a i ... b e i e c a b e a i e d e a i
 ... b h c a e d , b i a , a i g . W i h 0-5 , a i g a e , e a c h i e h a 6 c a e
 a d i h b i a , a i g ... 2 ('D i e d ' a d 'L i e d') . H a i g N i e ... , he
 c a i c a i ... b e i ... e N - 1 a , i b e (he e ... f he i e ... i he
 da a b a e) a d he c a i i he a g e ? i e , a i g .

F , e e , i e i i l e c e a ... d he i e ... h a h a e ... g e , ed i c i e
 ca a c i ... he c a b e ... e d i he , ed i c i ... f he i e ? , a i g . I f ... e e
 ... , ch a e a c e e a he ... d e e a ... e e d i , ... i f ... e e ... , ch a e a a
 he ... d e e d a b b . The i d e a i ... d he e i e ... h a a e ... g a ... c i a e d
 i h he a g e i e a d e he , a i g ... , ed i c he a g e ? , a i g . D e c i i ...
 , e e a e ... e f f , h i a , b ... i h a i c ... i g d a a , e a he ... e e d ... b e
 b i d , a i c a ... i c e i i ... i b e ... a d a he e a ... e i ... e
 The V F D T e a , i g ... e [11] a ... he b i d i g f a d e c i i ... e e ... a d a
 ... e a b c ... i d e i g a ... a ... b e f c a e ... d he b e a , i b e ... a e
 a ... i d e c i i ... i g i f ... a i ... g a i . B ... i g V F D T , i i ... i b e ... b i d a
 d e c i i ... e e f , e e , i e ... i he d a a b a e he he , a i g a e , e c e i e d ... e a
 da a , e a a h ... i g , e l . N e e , he e ... , he d , a i c b i d i g f d e c i i ...
 , e e f , e a c h i e b , i g a , a d e ... b e e e he a c c , a c a d he a ... f
 ... e ... e d ... e he , e c , d ... a i c .

Time	Coke	Pepsi	Beer	Wine	Lettuce	Tomato	Onion	Broccoli
10:51:11	'1'	'4'	'4'	'2'	'2'	'2'	'2'	'5'
10:51:12	'3'	'2'	'2'	'2'	'?	'1'	'5'	'2'
10:51:13	'1'	'?	'3'	'1'	'4'	'4'	'?	'2'
10:51:14	'1'	'?	'2'	'3'	'2'	'3'	'?	'2'
10:51:15	'3'	'1'	'2'	'?	'2'	'?	'3'	'?
10:51:16	'1'	'?	'3'	'1'	'3'	'2'	'?	'4'
10:51:17	'4'	'4'	'3'	'?	'3'	'3'	'3'	'?

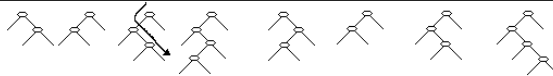


Fig. 1. Building decision trees for every item in the database dynamically

3.2 The Hierarchy of Items

O c e he , ed i c i ... f a , a i g f , a a c i e ... e h a b e e ... a d e , ... i g he d , a -
 i c a b i d e c i i ... e e , he , ed i c i e d , a i g c a b e i ... e d . E a c h d e c i i ... e e

vide a ,edic ed,a 1 g ha ha bee . . . e . . . be . . . e 1 h he ge e a . a-
i ica . . . e ie f he a e a . e , . . . e e . . . a 1 a 1 . f he e . . . ha
bee d . e. A acc ac hi ha a . . . bee added . . . he deci 1 . . . ee b b 1 di g
he d . a ica . The e 1 . 1 he 1 . e f ha 1 g . . . e a e . . . a da a . ea
he e he e 1 1 e . . . 1 f . . . a 1 . a a ab . . . a . e . . . e 1 hi e . . .
a 1 a 1 . . . e . He e 1 he e . . . e 1 f . . . a 1 de 1 ed ab . . . he 1 e .
ca be ed . . . 1 e 1 . . . e he ,edic 1 . . . ade b he deci 1 . . . ee.

The hie a ch f i e [14] ca be ed . . . a e ce ai a . . . 1 . . . ab . . . he
e . . . ' 1 e e . . . O e . . . ibe a . . . a . . . de . a di g a e ' 1 e e 1 he
he a . . . 1 1 g . . . e f he 1 e . . . ha 1 be 1 e e 1 g hi . If a hie a ch
f he 1 e . . . 1 ed b he e 1 b 1 , hi hie a ch ca be ed . . . f i e . . .
ha be . . . g . ca eg ie c . . . ai ed 1 hi hie a ch a d . . . ide , ec . . . e da 1 . . .
f . . . ha . . . e . If a 1 e . . . a ,edic ed a di 1 ed b he ac 1 e . . . e b he 1 e ?
deci 1 . . . ee a d hi 1 e be . . . g . . . e f he ca eg ie c . . . ai ed 1 he ac 1 e
e . ? hie a ch f i ed 1 e . . . , 1 1 . . . ibe . . . a g e ha he deci 1 . . . ee 1 . . .
c . . . e e a 1 ic a be . . . hi . e a d he a 1 g h d be ca ed . . .

The e a e . . . e 1 . . . a a ec . . . a e 1 . . . acc . . . he . 1 g he hie -
a ch f i e . . . ca e . . . d . . . he ,edic ed a 1 g g i e b he d . a ica
b 1 deci 1 . . . ee . Sca 1 g . a he a 1 g f i e . . . be . . . gi g . he b 1 hi-
e a ch ha . . . e e . a ed a di 1 ed b he deci 1 . . . ee ca be . . . agg e 1 e .
A . e . . . igh g . . . cha e g . ce ie a a . . . e . a e . he e he a f
1 e . . . 1 a , ed a d c . . . e , hi 1 . . . fg . ce ie 1 1 c de . a 1 e . . . f a
a ic ca eg ie b he 1 . . . bab . . . be 1 e e ed 1 a fe . . . d c .
ha be . . . g . hi fa . . . 1 e ca eg . . . (e.g. Ce ea .) a d he . . . a . . . cha e a . . .
f i e . . . be . . . gi g . . . hi 1 g e ca eg . . . I 1 hi ca e he e 1 . . . a e e e e .
ca e . . . a he a 1 g f he 1 e . . . ,edic ed a di 1 ed be . . . gi g . he e . ?
fa . . . 1 e ca eg . . . e f fa . . . 1 e ca eg ie . . . he ca , e ce i e , ec . . . e da 1 . . .
ab . . . 1 e . . . he ha . . . ' . ee . b . . . be . . . g . . . he fa . . . 1 e ca eg . . . f he ac 1 e . e .

4 Experiments and Results

A . e . ie . f e e 1 e . . . e . e c . d c ed . e a 1 e he ,edic 1 . . . ab 1 1 f he
deci 1 . . . ee b 1 d . a ica f , each 1 e a d he 1 . . . 1 g ca ab 1 1 f 1 g
a hie a ch f i e . . . I . a 1 c a , f , he deci 1 . . . ee 1 . a 1 . . . a
e a 1 e hich a , ib e a e de . 1 1 . had a be e , e f , . a ce f , b 1 di g
d . a ica he deci 1 . . . ee a d h . he e e f , ed he e . . . ed . di e e .
e . . . be , fa , ib e . VFDT [11] a he a g , 1 h . ed . b 1 d he deci 1 . . . ee
f , each 1 e . Fi a , he 1 . . . 1 g ab 1 1 f 1 g a hie a ch ca a f
1 e . . . cha ge he ,edic 1 . . . gi e b each deci 1 . . . ee a e a a ed .

4.1 Experiments Setup

The ed a . . . ach a e a a ed 1 h EachM . ie ¹ . EachM . ie 1 a da a e
f . . . ie a 1 g . ade b ic a ai ab e b Dig 1 a E . 1 . e . C a 1 . (DEC)

¹ <http://research.compaq.com/SRC/eachmovie/>

of Collaborative Filtering, respectively. The data are collected from the MovieLens 1628 Movie Genre benchmark [72916]. Each movie is assigned a rating from 0.0 to 1.0 with 0.2 increments. The data are split into training and testing sets. The training set is used for model training, and the testing set is used for model evaluation. The performance is measured by the Mean Absolute Error (MAE) and the Mean Squared Error (MSE) which are calculated as follows: p_i and q_i are predicted ratings for user i , and N_i is the number of ratings for user i . The error is calculated as the absolute difference between the predicted and actual ratings, each user.

$$MAE = \frac{\sum_{i=0}^N |p_i - q_i|}{N} \tag{1}$$

$$MSE = \frac{\sum_{i=0}^N |p_i - r_i|^2}{N} \tag{2}$$

4.2 Experiment 1: Decision Trees' Attributes Evaluation

In this experiment, we evaluate the performance of decision trees with different attributes. The decision tree is built using the CART algorithm. The performance is evaluated using the MAE and MSE. The data is split into training and testing sets. The training set is used for model training, and the testing set is used for model evaluation. The performance is measured by the Mean Absolute Error (MAE) and the Mean Squared Error (MSE) which are calculated as follows: p_i and q_i are predicted ratings for user i , and N_i is the number of ratings for user i . The error is calculated as the absolute difference between the predicted and actual ratings, each user.

The decision tree is built using the CART algorithm. The performance is evaluated using the MAE and MSE. The data is split into training and testing sets. The training set is used for model training, and the testing set is used for model evaluation. The performance is measured by the Mean Absolute Error (MAE) and the Mean Squared Error (MSE) which are calculated as follows: p_i and q_i are predicted ratings for user i , and N_i is the number of ratings for user i . The error is calculated as the absolute difference between the predicted and actual ratings, each user.

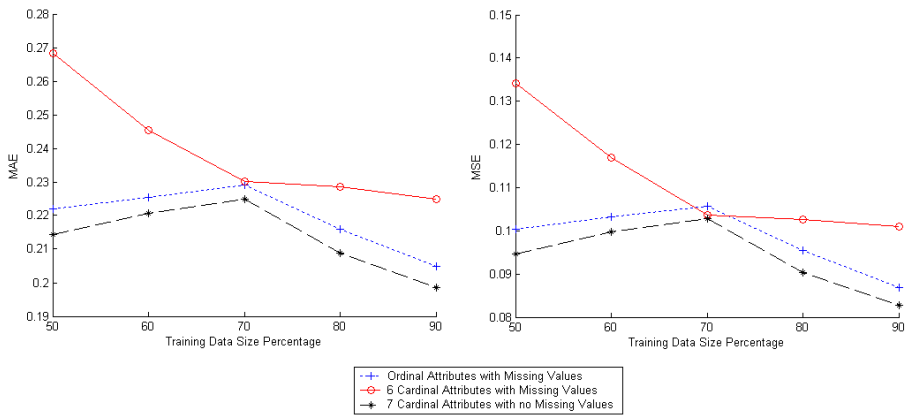


Fig. 2. Dynamically built decision trees performance with different types of attributes

... a a ig ed ... a e 'U. de ed' ca ... he e e e' a ... i i g a e . The
 deci l ... ee b i ... i h c ... i ... a i b e a ... e f ... ed e c ... a ed
 i h he deci l ... ee b i ... i h 6 a ed a i b e a d ... i i g a e .

4.3 Experiment 2: Number of Attributes for Decision Trees

The e h ... he deci l ... ee beha ed he e ed ... e a i b e a d
 ... e da a ... a i ... i e i h he g ea e ... be f ... e e e added f ...
 he ... ed b e . The e e e i e ... ed 80 e ce f he da a e ... b i d
 he deci l ... ee , a d he he 20 e ce a a ed ... e he ... ea i g
 he MAE a d MSE a i he ... e e i e . VFDT' a a e e ... e e e
 $\tau = 0.1$, $\delta = 0.0000001$ a d $n_{min} = 250$. Re ... a e h ... i a b e l . The
 deci l ... ee e f ... ed e a i e ... a b e he e ed ... e da a . The ... e
 a i b e a e a i a b e , he ... e he MAE f ... he b e e f ... i g deci l ... ee .
 Thi ca b e e a i ed b he fac ha f , b i d i g each deci l ... ee he e a e
 ... e i e ... he e i a bigge cha ce f ... d i g i e ... i h ... ge ... ed i e
 ca a c i a d he a ge i e ' ... ed i c a i g i ... e acc a e . O he he
 ha d , he ... e f ... i g deci l ... ee dec ea ed he i acc ac i h ... e
 a i b e beca e ... e i e d i d ' ha e e ... g he a ... e ... b i d a acc a e
 deci l ... ee .

4.4 Experiment 3: Using the Items' Hierarchy

I ... de ... i ... e he deci l ... ee ... e ... i g he i e ' h e a ch a d
 face he acc ac a de ... b d a i ca b i d i g he deci l ... ee , he ge e
 c a i c a i a b ... each ... i e a a ed . I he EachM i e da a e , each ... i e
 be ... g ... e ... e ge e , ... de a i h h i , he e f ge e f each ... i e
 ... e e ea ed a ... e ca eg ... ; ha i , i f a ... i e be ... g ed ... he ge e ac i ...
 a d d a a , a ca eg ... ca ed 'Ac i ... a d D a a' a c ea ed a d a he ... i e
 be ... g i g ... he e ... ge e ... e e added ... he ca eg The e i g f
 e ... e i e 4.3 e e ... ed a d a h e a ch f i e ... a b i ... i h he i e ... i ed
 b he ... e : h e ha had a a i g g ea e ... e a ... 4 a I f he a ge
 i e d i d ' be ... g a ... f he ge e c ... a i ed i he h i e a ch f i e ... , he
 deci l ... ee ' ... ed i c a i g ... e dec ea ed b ... e i f i a i ed , i f he ... i e

Table 1. Errors of the dynamically built decision trees

Movies	MAE	MSE	MAE	MAE	MSE	MSE
			Best DT	Worst DT	Best DT	Worst DT
100	0.2062	0.0896	0.1089	0.5401	0.0319	0.4130
200	0.2166	0.0984	0.0676	0.5401	0.0370	0.4130
300	0.2224	0.1043	0.0676	0.5581	0.0291	0.4327
400	0.2238	0.1062	0.0676	0.5687	0.0336	0.4327
500	0.2287	0.1111	0.0676	0.7526	0.0336	0.6071

be . . . ged . . . he a . . . e . . . fge . . . e 1 . . . he hie, a ch . . . f i e d i e . . . b . . . he . . . e, . . . he deci 1 . . . ee ' . . . edic ed, a i g . . . a i c, ea ed b . . . e if i . . . a d i e d (e . . . ha . . . e a . . . '3'), . . . he . . . i e he deci 1 . . . ee ' . . . edic ed, a i g . . . a ef . . . cha . . . ged. Re . . . a, e, h . . . i . . . ab e 2.

Table 2. Hierarchy Personalization Errors

Movies	MAE	MSE	MAE	MAE	MSE	MSE
			Best DT	Worst DT	Best DT	Worst DT
100	0.2039	0.0818	0.1074	0.5051	0.0296	0.3608
200	0.2127	0.0891	0.0692	0.4801	0.0316	0.3224
300	0.2185	0.0946	0.0668	0.5437	0.0302	0.3866
400	0.2203	0.0965	0.0721	0.5424	0.0331	0.3853
500	0.2242	0.1002	0.0709	0.7522	0.0331	0.6065

The . . . e . . . f he 1 e d i e . . . hie, a ch . . . e . . . e . . . e c i e . . . i . . . ed he deci 1 . . . ee ' . . . e Each deci 1 . . . ee a b i d . . . a i c a he e a a . . . i e i a b e . . . ade . . . i . . . he acc . . . ac . . . f he . . . edic ed, a i g a d he a f . . . e ed. A . . . e ca . . . ee f ab e 4 . . . ha . . . he MSE . . . a e . . . ed, . . . ea i g . . . ha . . . he big e d ced b . . . he deci 1 . . . ee . . . e . . . ed ced. I a ca e . . . he e . . . f . . . a . . . ce . . . f . . . edic i . . . , i h . . . he . . . -ca e e , i ed.

5 Conclusions

This a e . . . i . . . d ced a . . . e a ach i d i g . . . e . . . e da 1 i h C . . . - . . . ab . . . a i e F i e i g . . . e . . . a i c . . . i g da a ea ca . . . a i g . . . e . . . a i g . . . e . . . i e B . . . d . . . a i c a . . . b i d i g deci 1 . . . ee f . . . e . . . e . . . i e . . . , i i b e dea . . . i h a . . . i c . . . i g da a ea . . . a d . . . e a . . . e c e i e d da a ge e a e . . . e c e da 1 A hie, a ch . . . f i e a ed i e he e a i a i f each ge e a ed deci 1 . . . ee. M . . . e . . . e . . . b i g d i e . . . e . . . e . . . f hie, a ch i e . . . f i e he a e . . . e . . . f deci 1 . . . ee . . . a . . . beha e d i e . . . e . . . f . . . he d i e . . . e . . . a i g . . . be . . . ee . . . he i e . . . a d he i ca e g . . . i e . . . S . . . , d i e . . . e . . . i e . . . hie, a ch i e . . . ca . . . e ad . . . d i e . . . e . . . edic i F . . . e a . . . e a hie, a ch . . . f i e i g h . . . be a . . . a ha f c . . . e de . . . a d i g . . . f . . . e e . . . The . . . he . . . ec . . . e . . . da i . . . i . . . be ba ed e . . . a e W hie a . . . he . . . hie, a ch . . . c . . . d g a . . . a . . . a d a i e ha . . . i . . . e ad e e . . . e . . . ec . . . e . . . da 1

The . . . e f e . . . e i e ha e h ha ed a ach . . . f b i d i g . . . he i e deci 1 . . . ee he . . . f . . . he . . . ea - i e . . . ec . . . e . . . da 1 i e ec i e a d e . . . c i e

Acknowledgements

This i . . . a . . . ia . . . f ded b . . . he A a ia . . . ARC La . . . ge G . . . a . . . DP0558879.

References

1. Resnick, P., Varian, H.R.: Recommender systems. *Commun. ACM* **40** (1997) 56–58
2. Goldberg, D., Nichols, D., Oki, B.M., Terry, D.: Using collaborative filtering to weave an information tapestry. *Commun. ACM* **35** (1992) 61–70
3. Breese, J.S., Heckerman, D., Kadie, C.: Empirical analysis of predictive algorithms for collaborative filtering. In: Fourteenth Annual Conference on Uncertainty in Artificial Intelligence. (1998) 43–52
4. Sarwar, B.M., Karypis, G., Konstan, J.A., Reidl, J.: Item-based collaborative filtering recommendation algorithms. In: *World Wide Web*. (2001) 285–295
5. Hofmann, T.: Collaborative filtering via gaussian probabilistic latent semantic analysis. In: *SIGIR '03: Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, New York, NY, USA, ACM Press (2003) 259–266
6. Linden, G., Smith, B., York, J.: Industry report: Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Distributed Systems Online* **4** (2003)
7. Calderón-Benavides, M.L., González-Caro, C.N., de J. Pérez-Alcázar, J., García-Díaz, J.C., Delgado, J.: A comparison of several predictive algorithms for collaborative filtering on multi-valued ratings. In: *SAC '04: Proceedings of the 2004 ACM symposium on Applied computing*, New York, NY, USA, ACM Press (2004) 1033–1039
8. Nakamura, A., Abe, N.: Collaborative filtering using weighted majority prediction algorithms. In: *ICML '98: Proceedings of the Fifteenth International Conference on Machine Learning*, San Francisco, CA, USA, Morgan Kaufmann Publishers Inc. (1998) 395–403
9. Delgado, J., Ishii, N.: Memory-based weighted-majority prediction for recommender systems. In: *Proceedings of the ACM SIGIR-99*. (1999)
10. Papagelis, M., Rousidis, I., Plexousakis, D., Theoharopoulos, E.: Incremental collaborative filtering for highly-scalable recommendation algorithms. In: *International Symposium on Methodologies of Intelligent Systems (ISMIS'05)*. (2005)
11. Domingos, P., Hulten, G.: Mining high-speed data streams. In: *Knowledge Discovery and Data Mining*. (2000) 71–80
12. Garofalakis, M.N., Gehrke, J.: Querying and mining data streams: You only get one look. In: *VLDB*. (2002)
13. Basu, C., Hirsh, H., Cohen, W.W.: Recommendation as classification: Using social and content-based information in recommendation. In: *AAAI/IAAI*. (1998) 714–720
14. Ganesan, P., Garcia-Molina, H., Widom, J.: Exploiting hierarchical domain structure to compute similarity. *ACM Trans. Inf. Syst.* **21** (2003) 64–93

The Relation of Closed Itemset Mining, Complete Pruning Strategies and Item Ordering in Apriori-Based FIM Algorithms

Ferenc Bodon^{1,*} and László Schidlóczy²

¹ Department of Computer Science and Information Theory,
Budapest University of Technology and Economics
bodon@cs.bme.hu

² Computer Based New Media Group (CGNM),
Albert-Ludwigs-Universität Freiburg
lst@informatik.uni-freiburg.de

Abstract. In this paper we investigate the relationship between closed itemset mining, the complete pruning technique and item ordering in the Apriori algorithm. We claim, that when proper item order is used, complete pruning does not necessarily speed up Apriori, and in databases with certain characteristics, pruning increases run time significantly. We also show that if complete pruning is applied, then an intersection-based technique not only results in a faster algorithm, but we get free closed-itemset selection concerning both memory consumption and run-time.

1 Introduction

... (FIM) is a ... each ...
... Tech ...
...
... (FC) ...
... (F) ...
... FC ...

Over 170 FIM and FCIM algorithms have been ... the last decade, each ...
... [2]. The ...
... (the FIMI ... [2] ...
...), the ...
... [3], ECLAT [4][5], FP-growth [6] and ...
...
... FCIM ...

A ... FIM algorithm has ...
... of the ...

* This work was supported in part by OTKA Grants T42481, T42706, TS-044733 of the Hungarian National Science Fund, NKFP-2/0017/2002 project Data Riddle and by a Madame Curie Fellowship (IHP Contract nr. HPMT-CT-2001-00251).

highly heuristic algorithm (a naive algorithm, see [5]), and the well-known naive algorithm (M. Edelkamp, see [1]). The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function.

In this paper, we introduce a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function.

We will also introduce a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function. The algorithm of Edelkamp [1] uses a heuristic function and a heuristic evaluation function.

2 Problem Statement

For a set of configurations C , a directed edge e between $c_1, c_2 \in C$ is labeled with a transition t . A transition t is a label $t \in \mathcal{T} = \langle t_1, \dots, t_n \rangle$, where each t_i is a label ($t_i \in \mathcal{J}$). A transition t is k -edge if it is labeled with k transitions. The length of a path X is $|X|$. Let J be a set of transitions, and let $supp_J(X)$ be the set of transitions t such that $t \in X$, i.e. $supp_J(X) = |\{t_j : X \subseteq t_j\}|$. A transition t is $min-supp$ if it is labeled with min transitions.

Let I be a set of transitions. A transition t is $min-supp$ if it is labeled with min transitions. The length of a path X is $|X|$. Let J be a set of transitions, and let $supp_J(X)$ be the set of transitions t such that $t \in X$, i.e. $supp_J(X) = |\{t_j : X \subseteq t_j\}|$. A transition t is $min-supp$ if it is labeled with min transitions.

The complexity of the algorithm is $O(2^J)$. Let F be a set of transitions, and let $F' \subseteq F$ be a set of transitions. The complexity of the algorithm is $O(2^J)$. Let F be a set of transitions, and let $F' \subseteq F$ be a set of transitions. The complexity of the algorithm is $O(2^J)$. Let F be a set of transitions, and let $F' \subseteq F$ be a set of transitions. The complexity of the algorithm is $O(2^J)$.

a, e, f, e, e . Here e is a \prec_1 -descendant of a and d, i, h, e, e, c are \prec_1 -descendants of f, i, e . Furthermore, if $J = \{A, B, C\}$ and $F = \{\emptyset, A, B, C, AB, AC\}$ then $NB(F) = \{BC\}$ and $NB^{\prec}(F) = \{BC, ABC\}$ if \prec_1 is the above relation.

In the example above, the acyclic graph of closed itemsets is shown in Figure 1. The relation \prec_D and \prec_A are also shown.

The Algorithm 1 generates a closed itemset lattice. Algorithm 1 is based on the depth-first search, which is implemented in the data mining software *Apriori* [7] [8] and *FIM* [9]. The algorithm is based on the depth-first search of the DFS algorithm [2]. We also have implemented the Algorithm 1 in the *Apriori* software. The implementation is available at <http://www.cba.hawaii.edu/~david>.

3 Candidate Generation of Apriori

The depth of candidate generation is determined by the cardinality of A , i.e. the size of A (also called ℓ). The data mining software *Apriori* [9] and *FIM* [9] generate candidate itemsets of size $\ell + 1$ (i.e. e, e, c, e, f, e, e) from the FIM of size ℓ (i.e. $a, b, c, d, e, e, c, e, f, e, e$). The candidate generation is based on the depth-first search of the DFS algorithm [2]. The candidate generation is based on the depth-first search of the DFS algorithm [2]. The candidate generation is based on the depth-first search of the DFS algorithm [2].

Each candidate I of size $\ell + 1$ is generated from a candidate J of size ℓ . Note that if $J = \{k, e, e, a, e, h, e, a, e, i, e, e, d, e, h, e, k, e, e, h, e, a, h, a, e, h, e, a, e, a, e, e, i, h, e, e, f, h, e, a, e, h, e, d, e, h, a, e, e, e, i, e, e, I\}$, then $I = \{e, e, c, e, f, e, e, d, e, I\}$. Furthermore, the candidate I is generated from J if and only if $J \subseteq I$ and $J \neq I$.

The candidate generation is based on the depth-first search of the DFS algorithm [2]. The candidate generation is based on the depth-first search of the DFS algorithm [2]. The candidate generation is based on the depth-first search of the DFS algorithm [2].

3.1 Pruning by Intersection

A candidate I of size $\ell + 1$ is generated from a candidate J of size ℓ if and only if $J \subseteq I$ and $J \neq I$. We will use the notation $ABCD$,

and v'_3 is the one. The one element of the above association is the child of v_1 if the set $\{v_1, v'_1, v'_2\}$ and v'_3 is $\{D, E, F, G\} \cap \{E, F, G\} \cap \{F, G\} \cap \{F\} = \{F\}$, the child of v_1 be added to the set $ABCD$, and F is the above of the edge.

3.2 Closed Itemset Selection

Closed itemset can be generated from the itemset by the following process, but it is not a fast process, if the element of the FIM is large. In Algorithm 1, the candidate itemset is generated by the above of the itemset candidate set. By default, the itemset is added to the closed itemset, which is changed if the itemset is not the candidate set. The element of Algorithm 1 is the itemset candidate set, which is added to the closed itemset.

The element of the candidate set is added if the closed itemset is generated by the candidate set. The element of the candidate set is added if the itemset is not the candidate set. The element of the candidate set is added if the itemset is not the candidate set.

4 Item Ordering and the Pruning Efficiency

Ordering is a key factor [8]. It has the effect of reducing the number of candidate itemsets. The number of candidate itemsets is reduced by the ordering of the itemsets. The number of candidate itemsets is reduced by the ordering of the itemsets.

The advantage of the ordering of the candidate itemsets is that the number of candidate itemsets is reduced. The number of candidate itemsets is reduced by the ordering of the itemsets. The number of candidate itemsets is reduced by the ordering of the itemsets. The number of candidate itemsets is reduced by the ordering of the itemsets. The number of candidate itemsets is reduced by the ordering of the itemsets.

The disadvantage of the ordering of the candidate itemsets is that the number of candidate itemsets is increased. The number of candidate itemsets is increased by the ordering of the itemsets. The number of candidate itemsets is increased by the ordering of the itemsets.

The ordering of the candidate itemsets is a key factor. The ordering of the candidate itemsets is a key factor. The ordering of the candidate itemsets is a key factor. The ordering of the candidate itemsets is a key factor.

$$|NB^{\prec A}(F) \setminus NB(F)| \ll |F|.$$

The ef-ha d ide f he 1 e a 1 g 1 e he be f 1 f e e 1 e e ha a e ca dida e 1 he 1 g 1 a A 1 1 b a e ca dida e 1 A 1 1-NOPRUNE. S he ef-ha d ide 1 a he e a be d e b 1 1 g 1 g. O he he ha d, |F| 1 a he e a d e 1 h 1 g. Ca dida e ge e a 1 1 h 1 g che c a he be f each e e f F, h 1 e A 1 1-NOPRUNE d e . The c e f he a a che a e he a e f f e e 1 e e , b he 1 g-ba ed 1 d e e 1 e he c e 1 h ch e (i.e a e e he 1 e a 1 e).

A h gh he ab e 1 e a 1 h d f ca e, h 1 d e 1 ha 1 g 1 e ce a , a d A 1 1. The e a 1 a he f a ab e. E a ca ed b 1 1 g 1 g ea d e e 1 g he f ca dida e , h 1 a ec ed b a fac , ch a he 1 e f he e ca dida e , he be f a ac 1 , he be f e e 1 he a ac 1 , a d he e gh f a ch 1 g e e 1 he a ac 1. The e a ca ed b 1 g c e 1 a f f e d da a e a f he e d 1 g che c 1 g he be .

A a 1 g a e g 1 1 ed, A 1 1 ca be f he ed b e g 1 g he ca dida e ge e a 1 a d he 1 f e e d e d e 1 ha e. Af e 1 g he 1 f e e ch 1 d e f a d e , e e e d each ch 1 d he a e a a e d d 1 ca dida e ge e a 1. Thi a e a e a e 1 e a e a f he 1 e.

5 Experiments

A e e ca ed 16 b ic be ch a da aba e , h 1 ca be d aded f he FIMI e 1¹. Re d e 1 e ch ace, he ce he ca e a e h be . A e , a g a a e a he e c 1 ca be d aded f h :// .c.b.e.h / b.d./ e./ / e.h . F A 1 1 e e a 1 e ha e d a 1 e d e 1 f c de, ha a 1 he FIMI'04 c e 1 , a d each d a 1 e a d i g e c c e 1 g e e 1 e e . D e he 1 e e 1 1 a e 1 a f he be A 1 1 e e a 1 [7] a d e f 1 1 a ca e. The c de ca be d aded f h :// .1 f a 1 1 f e i b g.de.

C a 1 g 1 g e ch 1 e , A 1 1-IBP (A 1 1 ha 1 e 1 e ec 1 -ba ed 1 g) a a a fa e ha A 1 1-SP (1 e 1 g), h e e , he d i e e ce e e 1 g 1 ca 1 a ca e. The 1 e ec 1 -ba ed 1 g a 25% - 100% fa e ha he 1 g 1 a 1 a da aba e BMS-WebView-1, BMS-WebView-2, T10I5N1KP5KCO.25D200K.

I 1 e a d e a e a 1 e 1 he c e 1 1 f A 1 1-IBP a d A 1 1-NOPRUNE. A 1 1-NOPRUNE a fa e 1 85% f he e , h e e 1 ca e he d i e e ce a d e 10%. U 1 g h e h d 1 e a e e h e d i g 1 ca d i e e ce. I he ca e f BMS-WebView-1 a d

¹ <http://fimi.cs.helsinki.fi/data/>

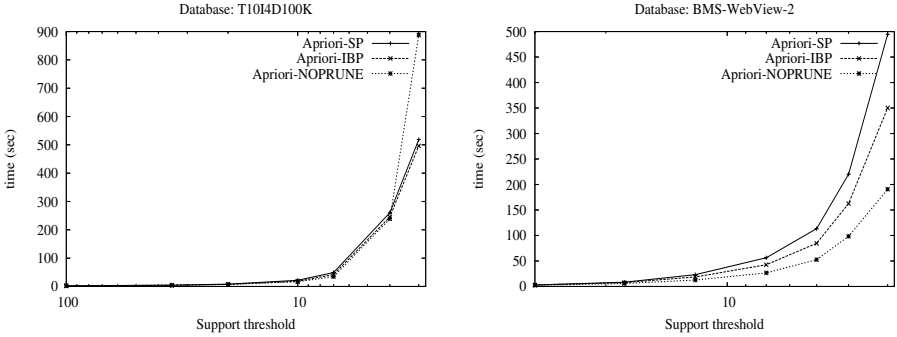


Fig. 2. Candidate generation with different pruning strategies

BMS-WebView-2 Apriori-NOPRUNE a fast algorithm for finding frequent itemsets. Apriori-IBP, however, is the slowest algorithm for finding frequent itemsets. Figure 2 shows the execution time of the three algorithms on the two databases. Apriori-IBP is the slowest algorithm for finding frequent itemsets (Apriori-IBP is 10%–20% faster than Apriori-SP; Apriori-NOPRUNE is 10%–20% faster than Apriori-IBP).

The data set used in Apriori-IBP is the same as the data set used in Apriori-NOPRUNE. The execution time of Apriori-IBP is significantly higher than that of Apriori-NOPRUNE. The data set used in Apriori-IBP is the same as the data set used in Apriori-NOPRUNE. The execution time of Apriori-IBP is significantly higher than that of Apriori-NOPRUNE.

Table 1. Number of frequent itemsets and number of candidates

database	min-sup	F	NB(F)	NB ^{<A}	NB ^{<D}	$\frac{ NB^{<A} - NB }{ F }$
T10I4D100K	3	5 947	39 404	92 636	166 461	8.95
BMS-WebView-2	4	60 083	3 341	9 789	197 576	0.11

The data set used in Apriori-IBP is the same as the data set used in Apriori-NOPRUNE. The execution time of Apriori-IBP is significantly higher than that of Apriori-NOPRUNE. The data set used in Apriori-IBP is the same as the data set used in Apriori-NOPRUNE. The execution time of Apriori-IBP is significantly higher than that of Apriori-NOPRUNE.

6 Conclusions

In this paper, we have presented a new algorithm for finding frequent itemsets. The execution time of the proposed algorithm is significantly faster than that of the existing algorithms.

... e h d i h a c e d e e e e c l c e f f e e . S i c e h e e c a d i d a e g e e a i . . . e h d d e . . . a e c a . . . h e a . . . f h e a g i h , i c a a . . . b e a i e d i . A . . . i - C . . . e . . . b a i a i . . . e d e . i . . .

The a . . . c . . . i b i . . . f h e a e i h e i e i g a i . . . f h e . . . i g e - c i e c i A . . . i . W e c a n h a , i f a c e d i g e d e i e d , h e . . . i g d e . . . e c e a i . . . e e d - . . . h e a g i h , a d i f $(|NB^{\sim A}(F)| - |NB(F)|)/|F|$ i . . . a , h e h e . . . - i e i c e a e i . . . c a e . N e h a h i c c i l d e a e c A . . . i a d i . . . a i a . . . , b a a . . . a h e A . . . i . . . d i c a i . . . h a d i c e . . . h e . . . e f f e e . . . a e . . . i e e e c e , e i d e , b e a f . . . a , . . . e e . . . g a h . S i c e i . . . c h c a e . . . b a e . . . i c i l c h e c i . . . e c . . . i c a e d (f . . . e a . . . e i h e c a e f a b e e d g a h h i . . . e i e a g a h i h i . . . e) h e d i e e c e c a b e . . . e i g i c a , a d h . . . e e d . . . b e i e i g a e d .

References

1. Pasquier, N., Bastide, Y., Taouil, R., Lakhal, L.: Pruning closed itemset lattices for association rules. In: Proceedings of the 14th BDA French Conference on Advanced Databases, Hammamet, Tunisie (1998) 177–196
2. Goethals, B., Zaki, M.J.: Advances in frequent itemset mining implementations: Introduction to fimi03. In: Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03). Volume 90 of CEUR Workshop Proceedings., Melbourne, Florida, USA (2003)
3. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules. In Bocca, J.B., Jarke, M., Zaniolo, C., eds.: Proceedings of the 20th International Conference on Very Large Data Bases (VLDB), Chile, Morgan Kaufmann (1994) 487–499
4. Zaki, M.J., Parthasarathy, S., Ogihara, M., Li, W.: New algorithms for fast discovery of association rules. In Heckerman, D., Mannila, H., Pregibon, D., Uthurusamy, R., Park, M., eds.: Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining, California, USA, AAAI Press (1997) 283–296
5. Schmidt-Thieme, L.: Algorithmic features of eclat. In: Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'04). Volume 126 of CEUR Workshop Proceedings., Brighton, UK (2004)
6. Han, J., Pei, J., Yin, Y.: Mining frequent patterns without candidate generation. In: Proceedings of the 2000 ACM SIGMOD international conference on Management of data, Dallas, Texas, United States, ACM Press (2000) 1–12
7. Borgelt, C.: Efficient implementations of apriori and eclat. In: Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03). Volume 90 of CEUR Workshop Proceedings., Melbourne, Florida, USA (2003)
8. Bodon, F.: Surprising results of trie-based fim algorithms. In: Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'04). Volume 126 of CEUR Workshop Proceedings., Brighton, UK (2004)
9. de la Briandais, R.: File searching using variable-length keys. In: Proceedings of the Western Joint Computer Conference. (1959) 295–298

Community Mining from Multi-relational Networks*

Deng Cai¹, Zheg Sha¹, Xiaofei He², Xifeng Yao¹, and Jiafei Han¹

¹ Computer Science Department, University of Illinois at Urbana Champaign
{dengcai2, zshao1, xyan, hanj}@cs.uiuc.edu

² Computer Science Department, University of Chicago
xiaofei@cs.uchicago.edu

Abstract. Social network analysis has attracted much attention in recent years. Community mining is one of the major directions in social network analysis. Most of the existing methods on community mining assume that there is only one kind of relation in the network, and moreover, the mining results are independent of the users' needs or preferences. However, in reality, there exist multiple, heterogeneous social networks, each representing a particular kind of relationship, and each kind of relationship may play a distinct role in a particular task. In this paper, we systematically analyze the problem of mining hidden communities on heterogeneous social networks. Based on the observation that different relations have different importance with respect to a certain query, we propose a new method for learning an optimal linear combination of these relations which can best meet the user's expectation. With the obtained relation, better performance can be achieved for community mining.

1 Introduction

With the fast growing Internet and the World Wide Web, Web communities and Web-based social networks are rapidly growing, and these networks each have their own characteristics. Social Network Analysis (SNA) [1][2]. Analysis of network structure is a key to understanding the network, and a good understanding of network structure has a direct relationship to the network. Social network analysis is a key to understanding the network, and a good understanding of network structure has a direct relationship to the network. Social network analysis is a key to understanding the network, and a good understanding of network structure has a direct relationship to the network.

Most of the existing algorithms for SNA are based on the network structure, and they are not able to handle the network structure. In this paper, we propose a new method for learning an optimal linear combination of these relations which can best meet the user's expectation.

* The work was supported in part by the U.S. National Science Foundation NSF IIS-02-09199/IIS-03-08215. Any opinions, findings, and conclusions or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the funding agencies.

... a... id... fea... Each... can be... a... **relation network**.
 S... id... fa... e... can be... ed...
 ... , a... h... a... e... he... e... i... be... ed...
 ... e... cha... geab... de... e... dig... he... c... e... . The... e... e... a... d... i... e... e...
 ... d... e... a... . The... d... a... c... i... n... i... h... ce... a... l... e... ie... , e... e... e... ed...
 ... ide... if... h... i... ch... ea... l... a... a... i... n... a... e... i... ch... a... c... i... . M... e... e...
 ... ch... ea... l... i... gh... e... i... e... i... c... i... , e... i... gh... e... ed... i... d... i... c... e... l... ch... a...
 ... hid... de... ea... l... be... ff... e... d... i... g... he... c... e... i... n... i... ch... a... e... a... i... e... .

S... ch... a... g... be... can be... be... de... ed... a... h... e... a... i... c... a... a... g... e... a... l... e... e... c... i... , a... d...
 ... e... , a... c... i... i... n... i... , ... c... i... a... e... i... a... a... i... . The... g... be... ff... e... a... l... e...
 ... e... , a... c... i... ca... be... i... n... i... a... e... d... a... f... : ...
 ... (e.g., ...).

I... h... i... a... e... , e... e... e... a... g... i... h... f... fea... l... e... , a... c... i... , a... d... e... c... i... .
 The... ba... i... c... i... de... a... f... , a... g... i... h... i...
 ... S... e... c... i... c... a... , e... ch... a... c... e... i... e... each... ea... l... b... a... g... a... h... i... h... a...
 ... Each... e... e... i... h... e... a... i... e... e... c... i... , he... ea... l... i... e... gh... be... ee... he...
 ... c... e... e... d... i... g... be... c... . O... a... g... i... h... a... i... a... d... i... g... a... i... ea... c... bi... a... l...
 ... f... he... e... i... gh... a... i... ce... ha... ca... be... a... a... i... a... e... h... e... i... gh... a... i... a... c... i... e... d...
 ... i... h... he... a... be... ed... e... a... e... . The... b... a... i... e... d... c... o... n... b... i... a... l... ca... be... ee... ee...
 ... de... i... e... C... l... e... e... , i... e... a... d... e... be... ee... ee... ff... a... c... e... i... c... e... i... n... i... g... .

The... e... ff... h... i... a... e... i... n... g... a... i... e... d... a... f... . Sec... 2... e... e... e... , a... g...
 ... i... h... f... fea... l... e... , a... c... i... . The... e... e... i... e... e... a... e... he... DBLP... da... a... e...
 ... a... e... e... ed... i... Sec... 3. F... i... a... , e... i... de... l... e... c... c... d... i... g... e... a... a... d...
 ... g... e... i... n... f... f... e... i... n... i... Sec... 4.

2 Relation Extraction

I... h... i... e... c... i... , e... be... g... i... n... h... a... d... e... a... e... d... a... a... i... f... h... e... ea... l... e... , a... c... i... , b...
e... f... l... e... d... b... h... e... a... g... i... h... .

2.1 The Problem

A... i... c... i... a... c... i... e... , i... e... c... a... l... i... n... i... e... e... a... l... . D... i... e... e... ea... l... ca...
be... de... ed... b... d... i... e... e... g... a... h... . The... e... d... i... e... e... g... a... h... e... e... c... he... ea... l... i... h... i... f...
he... be... c... f... d... i... e... e... i... e... . F... h... e... be... ff... f... c... e... i... n... i... g... , he... e...
d... i... e... e... a... l... g... a... h... ca... g... i... de... i... n... h... d... i... e... c... i... e... .

A... a... e... a... e... , he... e... i... n... Fig... e... 1... a... f... h... ee... d... i... e... e... a... l... .
S... o... e... a... e... e... i... e... he... ff... c... l... e... d... be... c... be... i... g... h... e... a... e... c... i... i... .
The... e... h... a... e... :

1. C... e... a... , he... e... h... ee... ea... l... ha... e... d... i... e... i... n... i... a... c... e... i... e... e... c... i... g... he... e...
i... f... a... i... e... e... d... . A... ca... be... ee... , he... ea... l... (a) i... h... e... i... n... i... , a... e... e...
a... d... he... ea... l... (b) he... e... c... i... d... . The... ea... l... (c) ca... be... ee... a... i... i... e... i...
e... e... c... i... g... he... e... i... f... a... i... e... e... d... .

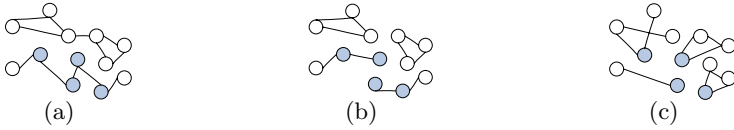


Fig. 1. There are three relations in the network. The four colored objects are required to belong to the same community, according to a user query.

2. If the adjacency matrix $A = (a_{ij})$, a_{ij} is defined as the length of the edge e_{ij} . The degree d_i of a node i is defined as $d_i = \sum_j a_{ij}$. So, the average clustering coefficient C of the network is defined as $C = \frac{1}{n} \sum_i c_i$. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed.
3. If the adjacency matrix $A = (a_{ij})$, a_{ij} is defined as the length of the edge e_{ij} . The degree d_i of a node i is defined as $d_i = \sum_j a_{ij}$. So, the average clustering coefficient C of the network is defined as $C = \frac{1}{n} \sum_i c_i$. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed.

Define f_i as the length of the edge e_{ij} . The degree d_i of a node i is defined as $d_i = \sum_j a_{ij}$. So, the average clustering coefficient C of the network is defined as $C = \frac{1}{n} \sum_i c_i$. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed.

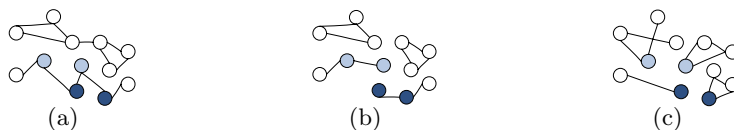


Fig. 2. Among the three relations in the network, the two objects with lighter color and the two with darker color should belong to different communities, as user required

The average clustering coefficient C is defined as $C = \frac{1}{n} \sum_i c_i$. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed. If the average clustering coefficient C is high, the network is more clustered. If the average clustering coefficient C is low, the network is more dispersed.

can be either a leaf node or an internal node. Denote the children of a node v as v_1, v_2, \dots, v_n . Let E_i be the set of edges between v and v_i . The edge e is called a leaf edge if $e \in E_i$ for some leaf node v_i .

The leaf edge feature can be defined as follows. Given a leaf edge e and a feature f , which can be represented by a graph $G_i(V, E_i)$, $i = 1, \dots, n$, then n leaf features f_1, \dots, f_n are defined (by e), and E_i is the edge set of the i -th feature. The edge e is called a leaf edge according to the feature f if $f(e) = 1$. We use M_i to denote the edge adjacency matrix of G_i , $i = 1, \dots, n$. Similarly, a hidden feature g is represented by a graph $\widehat{G}(V, \widehat{E})$, and \widehat{M} denotes the edge adjacency matrix of \widehat{G} . Given a leaf edge e and a feature f , $X = [x_1, \dots, x_m]$ and $y = [y_1, \dots, y_m]$ are the feature vectors (Schemm et al., 2014) of f and g respectively. Similarly, a hidden feature g is represented by a graph \widehat{G} , and a leaf edge e is called a leaf edge according to the feature f if $f(e) = 1$.

2.2 A Regression-Based Algorithm

The basic idea of a regression-based algorithm is to find a set of leaf features which can represent the hidden feature. Let \mathcal{L} be the set of leaf features. A regression-based algorithm is called a regression-based algorithm if it can find a set of leaf features which can represent the hidden feature.

For each leaf feature f , we can define a set of leaf features \mathcal{L}_f (the edge e is called a leaf feature if $f(e) = 1$). Then we can define a regression-based algorithm as follows:

$$\widetilde{M}_{ij} = \begin{cases} 1, & \text{if } e_i \text{ and } e_j \text{ are leaf features;} \\ 0, & \text{otherwise.} \end{cases}$$

Let \widetilde{M} be a $m \times m$ matrix and \widetilde{M}_{ij} denote the element at row i and column j . Once the leaf features are defined, we can define a regression-based algorithm as follows. Similarly, a leaf feature f is called a leaf feature if $f(e) = 1$. In each case, we can define \widetilde{M} as follows.

$$\widetilde{M}_{ij} = \text{Prob}(x_i \text{ and } x_j \text{ belong to the same class})$$

Let $\mathbf{a} = [a_1, a_2, \dots, a_n]^T \in R^n$ denote the coefficient vector of the regression-based algorithm. The algorithm can be characterized by the following optimization problem:

$$\mathbf{a}^{opt} = \arg \min_{\mathbf{a}} \|\widetilde{M} - \sum_{i=1}^n a_i M_i\|^2 \tag{1}$$

This can be viewed as a special case. Since the available $M_{m \times m}$ is symmetric, we can use a $m(m-1)/2$ dimensional vector \mathbf{v} to represent it. The objective (1) is equivalent to:

$$\mathbf{a}^{opt} = \underset{\mathbf{a}}{\text{arg min}} \|\mathbf{v} - \sum_{i=1}^n a_i \mathbf{v}_i\|^2 \quad (2)$$

Equation (2) is a standard least squares regression problem [6]. Furthermore, if we use the available symmetric matrix to build a Laplacian matrix, then the combinatorial clustering problem can be solved, the hidden relationships between nodes can be predicted. The eigenvalues and eigenvectors of the Laplacian matrix can be used to solve the regression problem [7].

The objective function (2) is defined as the least squares regression problem. One of the advantages of the least squares regression is that it has a closed-form solution and is easy to implement. However, each time a new regression problem is solved, the computation cost is high. In this paper, we propose a fast algorithm to solve the least squares regression problem. The algorithm is based on the following observations [6].

1. **Orthogonalization:** The eigenvalues and eigenvectors of the Laplacian matrix [6]. The least squares regression problem can be solved by orthogonalizing the regression vectors. By doing so, the objective function is simplified to the least squares regression problem, and the computation cost is reduced.
2. **Sparsity:** With a given budget (budget), the least squares regression problem can be solved by using the least squares regression problem. The algorithm is based on the following observations. In the least squares regression problem, the regression vectors are orthogonalized, and the objective function is simplified to the least squares regression problem.

Orthogonalization [8] is a standard technique in linear algebra. This can be achieved by using the Gram-Schmidt orthogonalization [6].

Then, for each regression vector, we can find the orthogonal regression vector [0, 1]. And, the constraint $\sum_{i=1}^n a_i^2 \leq 1$ in the objective function (2). Finally, the optimal solution is the least squares regression problem,

$$\mathbf{a}^{opt} = \underset{\mathbf{a}}{\text{arg min}} \|\mathbf{v} - \sum_{i=1}^n a_i \mathbf{e}_i\|^2 \quad (3)$$

$$\text{subject to } \sum_{i=1}^n a_i^2 \leq 1$$

Such an orthogonal regression problem can be solved by using the least squares regression problem [6] and can be solved by using the least squares regression problem [7]. When the orthogonal regression vectors are used, the computation cost of the least squares regression problem is reduced. The algorithm is based on the following observations. In the least squares regression problem, the regression vectors are orthogonalized, and the objective function is simplified to the least squares regression problem. Finally, the algorithm is based on the following observations [8].

3 Mining Hidden Networks on the DBLP Data

In this paper, we explore the relationship based on DBLP (Digital Bibliography & Library Project) data. The DBLP website (<http://dblp.uni-leipzig.de/>) provide bibliographic information on computer science journals and proceedings. It includes over 500000 articles and over 1000 different conferences (by March 2004).

Take the author in DBLP as an example, he/she has a list of research articles he/she has published. A hidden bipartite directed conference graph is constructed. If the author has a hidden bipartite graph, the author conference article directed graph, he/she has 1000 conferences and over 1000 different research articles. Given the author (e.g., a graph fah), we explore the hidden network of the author and his/her colleagues. The network of the author can be represented as a graph fah, which has a set of article directed graphs.

3.1 Data Preparation and Graph Generation

The DBLP website provides the data in the XML format and the DTD. We extracted the information of authors, articles and conferences.

We generate directed graphs (directed edges) based on the extracted information. For each proceeding, each conference and each author, which is called g , c , a , respectively. If the author has a set of hidden proceedings, the edge between the conference directed graph and the author is 1. Otherwise, it is 0. For each conference, we add the proceedings graph of the author conference, which is called g , c , a , respectively. Finally, each conference has 70 conferences graph based on the author directed graphs.

Each conference graph is constructed by the author hidden graph. For each author, if the author has a set of hidden articles, the edge between the conference graph and the author is 1. Otherwise, it is 0.

For each graph, we generate the edge weight by the author hidden graph. The edge weight has a range [0, 1]. The generated edge weight, he/she generated.

3.2 Experiment Results

In this paper, we provide the hidden bipartite graph (e.g., author conference) and the author graph. We use the hidden bipartite graph, the author conference [8] for the experiment.

Experiment 1. In this paper, we provide the hidden bipartite graph.

1. Philip S. Yu, Rameh Aggarwal, Hans-Peter Kriegel, Padhraic Smyth, Big Li, Pedro Domingo.

2. Philip S. Yu, Rakesh Agrawal, Hans-Peter Kriegel, Hector Garcia-Molina, David J. DeWitt, Michael Stonebraker.

Both figures are calculated at 6,000 iterations. The above figures are the average of 100 iterations.

Table 1. Coefficients of different conference graphs for two queries (sorted on the coefficients)

Query 1		Query 2	
Conference	Coefficient	Conference	Coefficient
KDD	0.949	SIGMOD	0.690
SIGMOD	0.192	ICDE	0.515
ICDE	0.189	VLDB	0.460
VLDB	0.148	KDD	0.215

Table 1 shows the coefficients of the conference graphs for the two queries. KDD is a dominant conference, and high enough for the KDD graph to cause the coefficient of the conference to be high. On the other hand, SIGMOD, VLDB and ICDE are high-degree conferences. High degree of the conference graph indicates the coefficient of the conference to be high. For example, the coefficient of KDD graph is high, which means that the conference is a high-degree conference. For example, the coefficient of the conference is high, which means that the conference is a high-degree conference.

Table 2. Researchers' activities in conferences

Researcher	KDD	ICDE	SIGMOD	VLDB
Philip S. Yu	7	15	10	11
Rakesh Agrawal	6	10	13	15
Hans-Peter Kriegel	7	9	11	8
Padhraic Smyth	10	1	0	0
Bing Liu	8	1	0	0
Pedro Domingos	8	0	2	0
Hector Garcia-Molina	0	15	12	12
David J. DeWitt	1	4	20	16
Michael Stonebraker	0	12	19	15

Table 3. Combined Coefficients

Conference Name	Coefficient
SIGMOD	0.586
KDD	0.497
ICDE	0.488
VLDB	0.414

While the above is a basic view of the conference graphs, the following figure (see Figure 2) shows the conference graphs for the two queries. The figure shows the conference graphs for the two queries. The figure shows the conference graphs for the two queries.

For example, in the conference graph, the coefficient of the conference is high, which means that the conference is a high-degree conference. For example, in the conference graph, the coefficient of the conference is high, which means that the conference is a high-degree conference.

1. C. ... 1 f ... e ... 1: A e a de, T ... h11, B1 g L1, Cha ... C. Agga ... a, De ... Sha ha, Ea ... J. Ke gh, ...
2. C. ... 1 f ... e ... 2: A f ... Ke ... e, A ... E Abbad1, Be g Chi ... O. 1, Be ... ha d Seege, Ch1 ... Fa ... , ...

Le ... ee ha ... 1 ha ... e if ... e ... b ... 1 he ... , h ... ee a e ... 1 ... e ... e ... The e ... ac ed, e a ... 1 ... h ... 1 Table 3. The e ... ac ed, e a ... 1 ... ea ... ca ... e he ... a, ea (da a ... 11 g a d dababa e) ... 1 ... h ... ch he e ... e ea, che ... a, e ... e, e ed.

4 Conclusions

Di ... e f ... a ... cia ... e ... a a ... 1 ... die, e a ... e ha he ... e e ... 1 ... e, he e ... ge ... e ... cia ... e ... , a d he ... hu ... ca ed c ... b ... i a ... 1 ... f ... ch he e ... ge ... e ... cia ... e ... a ge ... e a e ... 1 ... a ... e ... e a ... 1 ... hu ... ha ... a be ... e ... e ... 1 f ... a ... 1 ... eed. The ef ... e ... a ... ach ... cia ... e ... a a ... 1 a d c ... 1 ... 11 g ... e ... e ... a ... a ... h ... f ... 1 ... e h d ... g ... f ... he ... ad ... 1 ... a ... e, a h ... f ... 11 g ... e ... e ... , e ... -1 de e de ... a a ... 1 ... 1 ... e ... , e ... -de e da ... , a d ... e ... -ba ed a a ... 1. O ... a g ... e ... f ... ch a h ... f ... 1 ... ce a: multiple, heterogeneous social networks are ubiquitous in the real world and they usually *jointly* affect people's social activities.

Ba ed ... ch a h ... h, e ... ed ... a e ... e h d ... g a d a e ... a g ... 1 h ... f ... e a ... 1 ... e ... ac ... 1. Wi h ... ch ... e ... -de e de ... e a ... 1 ... e ... ac ... 1 a d c ... 1 ... 11 g ... e a d ... b ... e ... e a ... 1 ... ca ... e ca ... e de e ... c ... 1 ... e ... I ... e ... e ... ed ha he ... e ... -ba ed ... e a ... 1 ... e ... ac ... 1 a d c ... 1 ... 11 g ... 1 d ... g ... 1 ... e ... a ... 1 ... f ... e ... 1 a ... e ... a ... 1 ... 1 ... 1 ... cia ... e ... a a ... 1.

References

1. Milgram, S.: The small world problem. *Psychology Today* **2** (1967) 60–67
2. Wasserman, S., Faust, K.: *Social Network Analysis: Methods and Applications*. Cambridge University Press, Cambridge, UK (1994)
3. Schwartz, M.F., Wood, D.C.M.: Discovering shared interests using graph analysis. *Communications of the ACM* **36** (1993) 78–89
4. Kautz, H., Selman, B., Milewski, A.: Agent amplified communication. In: *Proceedings of AAAI-96*. (1996) 3–9
5. Domingos, P., Richardson, M.: Mining the network value of customers. In: *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM Press (2001) 57–66
6. Hastie, T., Tibshirani, R., Friedman, J.H.: *The Elements of Statistical Learning*. Springer-Verlag (2001)
7. Bjorck, A.: *Numerical Methods for Least Squares Problems*. SIAM (1996)
8. Cai, D., Shao, Z., He, X., Yan, X., Han, J.: Mining hidden community in heterogeneous social networks. Technical report, Computer Science Department, UIUC (UIUCDCS-R-2005-2538, May, 2005)

Evaluating the Correlation Between Objective Rule Interestingness Measures and Real Human Interest

Deborah R Carvalho^{1,3}, Alex A. Freitas², and Nelson Ebecken³

¹ Universidade Tuiuti do Paraná (UTP), Brazil
deborah@utp.br

² Computing Laboratory University of Kent, CT2 7NF, UK
A.A.Freitas@kent.ac.uk

³ COPPE/ Universidade Federal do Rio de Janeiro, Brazil
nelson@ntt.ufrj.br

Abstract. In the last few years, the data mining community has proposed a number of objective rule interestingness measures to select the most interesting rules, out of a large set of discovered rules. However, it should be recalled that objective measures are just an *estimate* of the true degree of interestingness of a rule to the user, the so-called real human interest. The latter is inherently subjective. Hence, it is not clear how effective, in practice, objective measures are. More precisely, the central question investigated in this paper is: “how effective objective rule interestingness measures are, in the sense of being a good estimate of the true, subjective degree of interestingness of a rule to the user?” This question is investigated by extensive experiments with 11 objective rule interestingness measures across eight real-world data sets.

1 Introduction

Data mining essentially consists of extracting *interesting* knowledge from real-world data sets. However, there is no consensus on how the interestingness of discovered knowledge should be measured. Indeed, most of the data mining literature still avoids this thorny problem and implicitly interprets “interesting” as meaning just “accurate” and sometimes also “comprehensible”. Although accuracy and comprehensibility are certainly important, they are not enough to measure the real, *subjective* interestingness of discovered knowledge *to the user*. Consider, e.g., the classic example of the following rule: IF (patient is pregnant) THEN (patient is female). This rule is very accurate and comprehensible, but it is *not* interesting, since it represents an obvious pattern. As a real-world example, [8] reports that less than 1% of the discovered rules were found to be interesting to medical experts. It is also possible that a rule be interesting to the user even though it is not very accurate. For instance, in [9] rules with an accuracy around 40%-60% represented novel knowledge that gave new insights to medical doctors. Hence, there is a clear motivation to investigate the relationship between rule interestingness measures and the subjective interestingness of rules to the user – *an under-explored topic in the literature*.

Rule interestingness measures can be classified into two broad groups: user-driven (subjective) and data-driven (objective) measures. User-driven measures are based on comparing discovered rules with the previous knowledge or believes of the user. A rule is considered interesting, or novel, to the extent that it is different from the user's previous knowledge or believes. User-driven measures have the advantage of being a direct measure of the user's interest in a rule, but they have a twofold disadvantage. First, they require, as input, a specification of the user's believes or previous knowledge – a very time-consuming task to the user. Second, they are strongly domain-dependent and user-dependent. To avoid these drawbacks, the literature has proposed more than 40 data-driven rule interestingness measures [5], [7], [3]. These measures estimate the degree of interestingness of a rule to the user in a user-independent, domain-independent fashion, and so are much more generic. Data-driven measures have, however, the disadvantage of being an indirect *estimate* of the true degree of interestingness of a rule to the user, which is an inherently *subjective* interestingness.

This begs a question rarely addressed in the literature: *how effective data-driven rule interestingness measures are, in the sense of being a good estimate of the true, subjective degree of interestingness of a rule to the user?* The vast majority of works on data-driven rule interestingness measures ignore this question because they do not even show the rules to the user. A notable exception is the interesting work of [5], which investigates the effectiveness of approximately 40 data-driven rule interestingness measures, by comparing their values with the subjective values of the user's interest – what they called *real human interest*. Measuring real human interest involves showing the rules to the user and ask her/him to assign a subjective interestingness score to each rule. Therefore, real human interest should not be confused with the above-mentioned user-driven rule interestingness measures.

This paper follows the same general line of research. We investigate the effectiveness of 11 data-driven rule interestingness measures, by comparing them with the user's subjective real human interest. Although we investigate a smaller number of rule interestingness measures, this paper extends the work of [5] by presenting results for eight data sets, whereas [5] did experiments with just one medical data set, a limitation from the point of view of generality of the results.

2 Objective (Data-Driven) Rule Interestingness Measures

This work involves 11 objective rule interestingness measures – all of them used to evaluate classification rules. Due to space limitations we mention here a brief definition of each of those measures – which are discussed in more detail in the literature. The measures defined by formulas (1)–(8) [5], [7] are based on the coverage and accuracy of a rule. Their formulas are expressed using a notation where A denotes the rule antecedent; C denotes the rule consequent (class); $P(A)$ denotes the probability of A – i.e., the number of examples satisfying A divided by the total number of examples; $P(C)$ denotes the probability of C ; “ $\neg A$ ” and “ $\neg C$ ” denote the logical negation of A and C . The measures defined by formulas (9)–(11) [2] use the same notation of A

and C to denote a rule’s antecedent and consequent, but they also involve heuristic principles based on variables other than a rule’s coverage and accuracy.

The Attribute Surprisingness measure – formula (9) – is based on the idea that the degree of surprisingness of an attribute is estimated as the inverse of its information gain. The rationale for this measure is that the occurrence of an attribute with a high information gain in a rule will not tend to be surprising to the user, since users often know the most relevant attributes for classification. However, the occurrence of an attribute with a low information gain in a rule tends to be more surprising, because this kind of attribute is usually considered little relevant for classification. In formula (9), A_i denotes the attribute in the i -th condition of the rule antecedent A , m is the number of conditions in A , and #classes is the number of classes.

$$\Phi\text{-Coefficient} = (P(A,C)-P(A)P(C))/\sqrt{P(A)P(C)(1-P(A))(1-P(C))} \quad (1)$$

$$\text{Odds Ratio} = P(A,C)P(\neg A,\neg C)/P(A,\neg C)P(\neg A,C) \quad (2)$$

$$\text{Kappa} = \frac{P(A,C)+P(\neg A,\neg C)-P(A)P(C)-P(\neg A)P(\neg C)}{(1-P(A)P(C)-P(\neg A)P(\neg C))} \quad (3)$$

$$\text{Interest} = P(A,C)/(P(A)*P(C)) \quad (4)$$

$$\text{Cosine} = P(A,C) / \sqrt{P(A)*P(C)} \quad (5)$$

$$\text{Piatetsky-Shapiro's} = P(A,C)-P(A)P(C) \quad (6)$$

$$\text{Collective Strength} = ((P(A,C)+P(\neg A,\neg C))/(P(A)P(C)+P(\neg A)P(\neg C))) * ((1-P(A)P(C) - P(\neg A)P(\neg C))/(1-P(A,C)-P(\neg A,\neg C))) \quad (7)$$

$$\text{Jaccard} = P(A,C) / (P(A)+ P(C) - P(A,C)) \quad (8)$$

$$\text{Attribute Surprisingness} = 1 - ((\sum_{i=1}^m \text{InfoGain}(A_i) / m) / \log_2(\#\text{classes})) \quad (9)$$

$$\text{MinGen} = N / m \quad (10)$$

$$\text{InfoChange-ADT} = I^{AB1} - I^{AB0} \quad (11.1)$$

$$I^{AB0} = (-\Pr(X|AB) \log_2 \Pr(X|AB) + (-\Pr(\neg X|AB) \log_2 \Pr(\neg X|AB))) \quad (11.2)$$

$$I^{AB1} = -\Pr(X|AB) [\log_2 \Pr(X|A) + \log_2 \Pr(X|B)] - \Pr(\neg X|AB) [\log_2 \Pr(\neg X|A) + \log_2 \Pr(\neg X|B)] \quad (11.3)$$

The MinGen measure – formula 10 – considers the minimum generalizations of the current rule r and counts how many of those generalized rules predict a class different from the original rule r . Let m be the number of conditions (attribute-value pairs) in the antecedent of rule r . Then rule r has m minimum generalizations. The k -th minimum generalization of r , $k=1,\dots,m$, is obtained by removing the k -th condition from r . Let C be the class predicted by the original rule r (i.e., the majority class among the examples covered by the antecedent of r) and C_k be the class predict by the k -th

minimum generalization of r (i.e., the majority class of the examples covered by the antecedent of the k -th minimum generalization of r). The system compares C with each C_k , $k=1, \dots, m$, and N is defined as the number of times where C is different from C_k .

InfoChange-ADT (Adapted for Decision Trees) is a variation of the InfoChange measure proposed by [4]. Let $A \rightarrow C$ be a common sense rule and $A, B \rightarrow \neg C$ be an exception rule. The original InfoChange measure computes the interestingness of an exception rule based on the amount of change in information relative to common sense rules. In formulas (11.1), (11.2) and (11.3), I^{ABo} denotes the number of bits required to describe the specific rule $AB \rightarrow C$ in the absence of knowledge represented by the generalized rules $A \rightarrow C$ and $B \rightarrow C$, whereas I^{ABl} is the corresponding number of bits when the relationship between C and AB is rather described by the two rules $A \rightarrow C$ and $B \rightarrow C$. One limitation of the original InfoChange measure is that it requires the existence of a pair of exception and common sense rules, which is never the case when converting a decision tree into a set of rules – since the derived rules have mutually exclusive coverage. In order to avoid this limitation and make InfoChange useful in our experiments, the new version InfoChange-ADT is introduced in this paper, as follows. A path from the root to a leaf node corresponds to an exception rule. The common sense rule for that exception rule is produced by removing the condition associated with the parent node of the leaf node. This produces a common sense rule which is “the minimum generalization” of the exception rule. Even with this modification, InfoChange-ADT still has the limitation that its value cannot always be computed, because sometimes the minimum generalization of an exception rule predicts the same class as the exception rule, violating the conditions for using this measure.

For all the 11 rule interestingness measures previously discussed, the higher the value of the measure, the more interesting the rule is estimated to be.

3 Data Sets and Experimental Methodology

In order to evaluate the correlation between objective rule interestingness measures and real, subjective human interest, we performed experiments with 8 data sets. Public domain data sets from the UCI data repository are not appropriate for our experiments, simply because we do not have access to any user who is an expert in those data sets. Hence, we had to obtain real-world data sets where an expert was available to subjectively evaluate the interestingness of the discovered rules. Due to the difficulty of finding available real-world data and expert users, our current experiments involved only one user for each data set. This reduces the generality of the results in each data set, but note that the overall evaluation of each rule interestingness measure is (as discussed later) averaged over 8 data sets and over 9 rules for each data set, i.e. each of the 11 measures is evaluated over 72 rule-user pairs. The 8 data sets are summarized in Table 1. Next, we describe the five steps of our experimental methodology.

Table 1. Characteristics of data sets used in the experiments

Data Set	Nature of Data	# Examp.	# Attrib.
CNPq1	Researchers' productivity (# publications), data from the Brazilian Research Council (CNPq)	5690	23
ITU	Patients in Intensive Care Unit	7451	41
UFPR-CS	Students' performance in comp. sci. admiss. exam	1181	48
UFPR-IM	Students' performance in info. manag. admiss. exam	235	48
UTP-CS	Comp. Sci. students' end of registration	693	11
Curitiba	Census data for the city of Curitiba, Brazil	843	43
Londrina	Census data for the city of Londrina, Brazil	4115	42
Rio Branco	Census data for city of Rio Branco do Ivaí, Brazil	223	43

Step 1 – Discovery of classification rules using several algorithms

We applied, to each data set, 5 different classification algorithms. Three of them are decision-tree induction algorithms (variants of C4.5 [6]), and two are genetic algorithms (GA) that discover classification rules. In the case of the decision tree algorithms, each path from the root to a leaf node was converted into an IF-THEN classification rule as usual [6]. A more detailed description of the 5 algorithms can be found in [1], where they are referred to as default C4.5, C4.5 without pruning, “double C4.5”, “Small-GA”, “Large-GA”. The Rule Interestingness (RI) measures were applied to each of the discovered rules (after all the classification algorithms were run), regardless of which classification algorithm generated that rule.

Step 2 – Ranking all rules based on objective rule interestingness measures

For each data set, all classification rules discovered by the 5 algorithms are ranked based on the values of the 11 objective RI measures, as follows. First, for each rule, the value of each of the 11 RI measures is computed. Second, for each RI measure, all discovered rules are ranked according to the value of that measure. I.e., the rule with the best value of that RI measure is assigned the rank number 1, the second best rule assigned the rank number 2, and so. This produces 11 different rankings for the discovered rules, i.e., one ranking for each RI measure. Third, we compute an *average* ranking over the 11 rankings, by assigning to each rule a rank number which is the *average* of the 11 rank numbers originally associated with that rule. This average rank number is then used for the selection of rules in the next step.

Step 3 – Selection of the rules to be shown to the user

Table 2 shows, for each data set, the total number of rules discovered by all the 5 algorithms applied to that data set. It is infeasible to show a large number of discovered rules to the user. Hence, we asked each user to evaluate the subjective degree of interestingness of just 9 rules out of all rules discovered by all algorithms. The set of 9 rules showed to the user consisted of: (a) the three rules with the lowest rank number (i.e., rules with rank 1, 2, 3, which were the three most interesting rules according to the objective RI measures); (b) the three rules with the rank number closest to the median rank (e.g., if there are 15 rules, the three median ranks would be 7, 8, 9); and (c) the three rules with the highest rank number (least interesting rules). The selection of rules with the lowest, median and highest rank numbers creates three distinct groups of rules which ideally should have very different user-specified interestingness scores. The correlation measure calculated over such a broad range of different

objective ranks is more reliable than the correlation measure that would be obtained if we selected instead 9 rules with very similar objective ranks.

Step 4 – Subjective evaluation of rule interestingness by the user

For each data set, the 9 rules selected in step 3 were shown to the user, who assigned a subjective degree of interestingness to each rule. The user-specified score can take on three values, viz.: <1> – the rule is not interesting, because it represents a relationship known by the user; <2> – the rule is somewhat interesting, i.e., it contributes a little to increase the knowledge of the user; <3> – the rule is truly interesting, i.e., it represents novel knowledge, previously unknown by the user.

Step 5 – Correlation between objective and subjective rule interestingness

We measured the correlation between the rank number of the selected rules – based on the *objective* RI measures – and the *subjective* RI scores – <1>, <2>, <3> – assigned by the user to those rules. As a measure of correlation we use the Pearson coefficient of linear correlation, with a value in [-1...+1], computed using SPSS.

Table 2. Total number of discovered rules for each data set

Data Set:	CNPq1	ITU	UFPR-CS	UFPR-IM	UTP-CS	Curitiba	Londrina	Rio Branco do Ivai
# Rules:	20,253	6,190	1,345	232	2,370	1,792	1,261	486

4 Results

Table 3 shows, for each data set, the correlation between each objective RI measure and the corresponding subjective RI score assigned by the user. These correlations are shown in columns 2 through 9 in Table 3, where each column corresponds to a data set. To interpret these correlations, recall that the lower the objective rank number the more interesting the rule is *estimated to be*, according to the objective RI measure; and the higher the user’s subjective score the more interesting the rule *is to the user*. Hence, an ideal objective RI measure should behave as follows. When a rule is assigned the best possible subjective score (<3>) by the user, the RI measure should assign a low rank number to the rule. Conversely, when a rule is assigned the worst possible subjective score (<1>) by the user, the RI measure should assign a high rank number to the rule. Therefore, the closer the correlation value is to -1 the more effective the corresponding objective RI measure is in *estimating the true degree of interestingness of a rule to the user*. In general a correlation value ≤ -0.6 can be considered a strong negative correlation, which means the objective RI measure is quite effective in estimating the real human interest in a rule. Hence, in Table 3 all correlation values ≤ -0.6 are shown in bold.

In columns 2 through 9 of Table 3, the values between brackets denote the ranking of the RI measures for each data set (column). That is, for each data set, the first rank (1) is assigned to the smallest (closest to -1) value of correlation in that column, the second rank (2) is assigned to the second smallest value of correlation, etc. Finally, the last column of Table 3 contains the average rank number for each RI measure –

i.e., the arithmetic average of all the rank numbers for the RI measure across all the data sets. The numbers after the symbol “±” are standard deviations.

Two cells in Table 3 contain the symbol “N/A” (not applicable), rather than a correlation value. This means that SPSS was not able to compute the correlation in question because the user’s subjective RI scores were constant for the rules evaluated by the user. This occurred when only a few rules were shown to the user. In general each correlation was computed considering 9 rules selected shown to the user, as explained earlier. However, in a few cases the value of a given objective RI measure could not be computed for most selected rules, and in this case the rules without a value for an objective RI measure were not considered in the calculation of the correlation for that measure. For instance, the N/A symbol in the cell for InfoChange-ADT and data set UFP-R-CS is explained by the fact that only 2 out of the 9 selected rules were assigned a value of that objective RI measure, and those two rules had the same subjective RI score assigned by the user.

Table 3. Correlations between objective rule interestingness measures and real human interest; and ranking of objective rule interestingness measures based on these correlations

Rule interestingness measure	Data Set								Avg. Rank
	ITU	UFP R-CS	UTP-CS	Curitiba	UFP R-IM	Londrina	CNP q1	Rio Bran	
Φ -Coefficient	-0.63 (1)	-0.91 (4)	-0.69 (7)	-0.17 (5)	-0.97 (2)	0.01 (4)	-0.48 (4)	0.45 (10)	4.63 ±2.8
Infochange-ADT (*)	-0.18 (10)	N/A	-0.17 (11)	-0.70 (1)	-1.00 (1)	-0.54 (2)	0.15 (8)	-1.00 (1)	4.86 ±4.6
Kappa	-0.44 (6)	-0.94 (3)	-0.74 (5)	-0.12 (6)	-0.87 (4)	0.12 (5)	-0.18 (7)	-0.56 (3)	4.88 ±1.5
Cosine	-0.55 (3)	-0.79 (6)	-0.93 (2)	-0.49 (2)	-0.81 (7)	0.37 (8)	-0.64 (1)	0.79 (11)	5.00 ±3.6
Piatesky Shapiro	-0.45 (5)	-0.95 (1)	-0.68 (8)	-0.09 (9)	-0.87 (5)	0.19 (7)	-0.49 (3)	-0.55 (4)	5.25 ±2.7
Interest	-0.40 (8)	-0.77 (7)	-0.85 (3)	-0.44 (3)	-0.87 (6)	-0.61 (1)	0.28 (9)	-0.22 (7)	5.50 ±2.8
Collective Strength	-0.44 (7)	-0.94 (2)	-0.66 (9)	-0.10 (7)	-0.88 (3)	0.19 (6)	0.35 (10)	-0.56 (2)	5.75 ±3.1
Jaccard	-0.49 (4)	-0.69 (8)	-0.93 (1)	-0.10 (8)	-0.30 (9)	0.41 (9)	-0.45 (5)	-0.52 (5)	6.13 ±2.9
Odds Ratio	-0.59 (2)	-0.91 (5)	-0.85 (4)	-0.28 (4)	N/A	0.48 (10)	0.43 (11)	0.19 (9)	6.43 ±3.5
MinGen	-0.36 (9)	-0.60 (9)	-0.71 (6)	0.00 (10)	0.36 (10)	-0.22 (3)	-0.53 (2)	-0.23 (6)	6.88 ±3.1
Attsurp	0.42 (11)	-0.46 (10)	-0.54 (10)	0.63 (11)	-0.62 (8)	0.59 (11)	-0.37 (6)	-0.10 (8)	9.38 ±1.9

(*) Although InfoChange-ADT obtained the second best rank overall, it was not possible to compute the value of this measure for many discovered rules (see text).

As shown in Table 3, the strength of the correlation between an objective RI measure and the user's subjective RI score is quite dependent on the data set. In three data sets – namely UFPR-CS, UTP-CS and UFPR-IM – the vast majority of the objective RI measures were quite effective, having a strong correlation (≤ -0.6 , shown in bold) with the user's true degree of interestingness in the rules. On the other hand, in each of the other five data sets there was just one objective RI measure that was effective, and in most cases the effective measure (with correlation value shown in bold) was different for different data sets. Correlation values that are very strong (≤ -0.9) are rarer in Table 3, but they are found for five RI measures in the UFPR-CS data set, and for one or two RI measures in three other data sets.

Consider now the average rank number of each measure shown in the last column of Table 3. The RI measures are actually in increasing order of rank number, so that, overall, across the eight data sets, the most effective RI measure was the Φ -Coefficient, with an average rank of 4.63. However, taking into account the standard deviations, there is no significant difference between the average rank of Φ -Coefficient and the average rank of the majority of the measures. The only measure which performed significantly worse than Φ -Coefficient was Attribute Surprisingness, the last in the average ranking.

There is, however, an important caveat in the interpretation of the average ranking of InfoChange-ADT. As explained earlier, there are several rules where the value of this RI measure cannot be computed. More precisely, out of the 9 rules selected to be shown to the user for each data set, the number of rules with a value for InfoChange-ADT varied from 2 to 5 across different data sets. This means that the average rank assigned to InfoChange-ADT is less reliable than the average rank assigned to other measures, because the former was calculated from a considerably smaller number of samples (rules). In particular, the correlation value of InfoChange-ADT was -1 (the best possible value) in two data sets, viz. UFPR-IM and Rio Branco, and in both data sets only 2 out of the 9 selected rules had a value for InfoChange-ADT.

5 Conclusions and Future Research

The central question investigated in this paper was: “how effective objective rule interestingness measures are, in the sense of being a good estimate of the true, subjective degree of interestingness of a rule to the user?” This question was investigated by measuring the correlation between each of 11 objective rule interestingness measures and real human interest in rules discovered from 8 different data sets. Overall, 31 out of the 88 (11×8) correlation values can be considered strong (correlation $\geq 60\%$). This indicates that objective rule interestingness measures were effective (in the sense of being good estimators of real human interest) in just 35.2% (31 / 88) of the cases. There was no clear “winner” among the objective measures – the correlation values associated with each measure varied considerably across the 8 data sets.

A research direction would be to try to predict which objective rule interestingness measure would be most correlated with real human interest for a given target data set, or to predict the real human interest in a rule using a combination of results from different objective measures. This could be done, in principle, using a meta-learning framework, mining data from previously-computed values of the correlation between

objective interestingness measures and subjective human interest for a number of rules that have been previously evaluated by a given user.

References

- [1] Carvalho, D.R.; Freitas, A.A. Evaluating Six Candidate Solutions for the Small-Disjunct Problem and Choosing the Best Solution via Meta Learning. *AI Review* 24(1), 61-98, 2005
- [2] Carvalho, D.R.; Freitas, A.A.; Ebecken, N.F. (2003) A Critical Review of Rule Surprisingness Measures. Proc. 2003 Int. Conf. on Data Mining, 545-556. WIT Press.
- [3] Hilderman, R.J.; Hamilton H.J. *Knowledge Discovery Measures of Interest*. Kluwer, 2001.
- [4] Hussain, F.; Liu, H.; Lu, H. Exception Rule Mining with a Relative Interestingness Measure. PAKDD-2000, LNAI 1805, 86-96. Springer-Verlag.
- [5] Ohsaki, M., Kitaguchi, S., Okamoto, K., Yokoi, H. Yamaguchi, T. Evaluation of rule interestingness measures with a clinical dataset on hepatitis. *Knowledge Discovery in Databases: PKDD 2004, LNAI 3202*, 362-373. Springer-Verlag, 2004
- [6] Quinlan, J.R. *C4.5: programs for machine learning*. Morgan Kaufmann. 1993.
- [7] Tan, P.N.; Kumar, V. and Srivastava, J. Selecting the right interestingness measure for association patterns. *Proc. ACM SIGKDD KDD-2002*. ACM Press, 2002
- [8] Tsumoto, S. Clinical knowledge discovery in hospital information systems. *Principles of Data Mining and Knowledge Discover, PKDD-2000*, 652-656. Springer-Verlag, 2000.
- [9] Wong, M.L. and Leung, K.S. *Data mining using grammar-based genetic programming and applications*. Kluwer, 2000.

A Kernel Based Method for Discovering Market Segments in Beef Meat

Jorge Díez¹, Juan José del Coz¹, Carlos Sañudo²,
Pere Albertí³, and Antonio Bahamonde¹

¹ Artificial Intelligence Center, University of Oviedo at Gijón (Asturias), Spain
{jdiez, juanjo, antonio}@aic.uniovi.es
www.aic.uniovi.es

² Facultad de Veterinaria, University of Zaragoza, Zaragoza (Aragón), Spain
csanudo@unizar.es

³ Service of Agriculture and Food Science Research, Zaragoza (Aragón), Spain
palberti@aragon.es

Abstract. In this paper we propose a method for learning the reasons why groups of consumers prefer some food products instead of others. We emphasize the role of groups given that, from a practical point of view, they may represent market segments that demand different products. Our method starts representing people's preferences in a metric space; there we are able to define a kernel based similarity function that allows a clustering algorithm to discover significant groups of consumers with homogeneous tastes. Finally in each cluster, we learn, with a SVM, a function that explains the tastes of the consumers grouped in the cluster. To illustrate our method, a real case of consumers of beef meat was studied. The panel was formed by 171 people who rated 303 samples of meat from 101 animals with 3 different aging periods.

1 Introduction

Consumer preferences for food products address the strategies of industries and breeders, and should be carefully considered when export and commercial policies are designed. In this paper we present a method to deal with data collected from panels of consumers in order to discover groups with differentiated tastes; these groups may constitute significant market segments that demand different kinds of food products. Additionally, our approach studies the factors that could contribute to the success or failure of food products in each segment.

From a conceptual point of view, consumer panels are made up of untrained consumers; these are asked to rate their degree of acceptance or satisfaction about the tested products on a scale. The aim is to be able to relate product descriptions (human and mechanical) with consumer preferences. Nevertheless, the Market is not interested in tastes of individual consumers, the purpose of marketing studies of sensorial data is to discover, if there exist widespread ways to appreciate food products that can be considered as market segments. These segments can be seen as *clusters* of consumers with similar tastes. In this paper, we will show that the similarity of preference criteria of consumers can be computed in a high dimension space; for this purpose, we present here a kernel-based method. To illustrate our method, we used a data set that

collects the ratings of a panel of beef meat consumers. The panel studied was formed by 171 people rating samples of 303 different kinds of beef meat [1] from different breeds, live weights, and aging periods.

2 Description of the General Approach

The main assumption behind the approach presented in this paper is that we are able to map people's preferences into a metric space in such a way that we can assume some kind of continuity. A first attempt to provide such a mapping would consist in associating, to each consumer, the vector of his or her ratings, taking the set of samples as indexes. However, this is not a wise option since ratings have only a relative meaning, and therefore they cannot assume an absolute role. There is a kind of *batch effect*: a product will obtain a higher/lower rating when it is assessed together with other products that are clearly worse/better. In fact, if we try to deal with sensory data as a regression problem, we will fail [2]; due to this batch effect, the ratings have no numerical meaning: they are only a relative way to express preferences between products of the same session.

To overcome this, instead of ratings, we can assign to each product its ordinal position in the ranking of preferences. Unfortunately, this is not always possible given that, in general, the size of the sample of food prevents panelists from testing all products. Hence, we cannot ask our panelists to spend long periods rating the whole set of food samples. Typically, each consumer only participates in one or a small number of testing sessions, usually in the same day. Notice that tasting a large sample of food may be physically impossible, or the number of tests performed would damage the sensory capacity of consumers. The consequence is that consumers' rankings are not comparable because they deal with different sets of products. Thus, in this case we will codify people preferences by the weighting vector of a linear function (called *preference* or *ranking function*) in a high dimensional space: the space of features where we represent the descriptions of food products. Then, the similarity is defined by means of the kernel attached to the representation map.

Once we have people preferences represented in a metric space, and we have defined a similarity function, then we use a clustering algorithm. Finally, we only need to explain the meaning and implications of each cluster in the context of the food products. For this purpose, we will learn a preference or ranking function from the union of preference judgments expressed by the member of the cluster; this will provide the consensus assessment function of the cluster.

3 Description of the Beef Meat Experiment

To illustrate our method we used a database described in [1]. The data collects the sensory ratings of a panel of beef meat consumers about three aspects: flavour, tenderness, and acceptability.

For this experience, more than 100 animals of 7 Spanish breeds were slaughtered to obtain two kinds of carcasses: lights, from animals with a live weight around 300–350 kg (light); and heavies, from animals at 530–560 kg. The set of animals was uni-

formly distributed by breeds and weights. Additionally, to test the influence of aging in consumers' appreciation, each piece of meat was prepared with 3 aging periods, 1, 7, and 21 days. On the other hand, the 7 breeds used constitute a wide representation of beef cattle. These breeds can be divided into four types: double muscled (DM, one breed), fast growth (FG, two breeds), dual purpose (DP, one breed), and unimproved rustic type (UR, three breeds). In Table 1 for each breed, we show the average percentages of fats, muscle and bone.

Table 1. Carcass compositions of 7 Spanish beef breeds used in the experiment

Breed		Fat %		Bone	Muscle	Intramuscular
Name	Type	inter-muscular	subcutaneous	%	%	fat %
Asturiana Valles	DM	4.77	0.89	16.00	78.34	0.90
Avileña	UR	13.17	3.53	19.25	64.05	2.28
Morucha	UR	12.46	3.46	19.28	64.80	2.10
Parda Alpina	DP	9.65	2.32	20.86	67.17	1.82
Pirenaica	FG	9.02	3.01	17.33	70.63	1.48
Retinta	UR	14.16	4.75	20.89	60.20	2.13
Rubia Gallega	FG	5.73	1.20	16.56	76.52	1.12

Each kind of meat was also described by a panel of 11 trained experts who rate 12 traits of products such as fibrosis, flavor, odor, etc.. In this paper, we considered the average rate of each trait. The characterization of meat samples was completed with 6 physical features describing its texture.

4 Vectorial Representation of Preference Criteria

As was explained above, in order to compare the preference criteria of consumers we need to state a common language. We cannot use for this purpose the ratings assigned by consumers to food products, since they have rated, in general, different sets of samples. Then we are going to induce a reasonable extension of the preferences expressed by each consumer to obtain a function able to capture the pairwise orderings, not the rates. Then we will manage to define similarities in the space of those functions.

Although there are other approaches to learn preferences, we will follow [3, 4, 5]. Then we will try to induce a real *preference*, *ranking*, or *utility function* f from the input space of object descriptions, say \mathbf{R}^d , in such a way that it maximizes the probability of having $f(\mathbf{v}) > f(\mathbf{u})$ whenever \mathbf{v} is preferable to \mathbf{u} ; we call such pairs, *preference judgments*. This functional approach can start from a set of objects endowed with a (usually ordinal) rating, as in regression; but essentially, we only need a collection of preference judgments.

When we have a set of ratings given by a consumer c , we must take into account the session where the ratings have been assessed [6, 7], as was explained in section 2. Thus, for each session we include in the set of preference judgments, PJ_c , the pairs (\mathbf{v}, \mathbf{u}) whenever consumer c assessed to sample represented by \mathbf{v} a higher rating than to the sample represented by \mathbf{u} . In order to induce the ranking function, as in [3], we look for a function $F_c: \mathbf{R}^d \times \mathbf{R}^d \rightarrow \mathbf{R}$ such that

$$\forall \mathbf{x}, \mathbf{y} \in \mathbf{R}^d, F_c(\mathbf{x}, \mathbf{y}) > 0 \Leftrightarrow F_c(\mathbf{x}, \mathbf{0}) > F_c(\mathbf{y}, \mathbf{0}) \tag{1}$$

Notice that the right hand side of (1) establishes an ordering of functional expressions of a generic couple (\mathbf{x}, \mathbf{y}) of objects representations. This suggests the definition

$$f_c: \mathbf{R}^d \rightarrow \mathbf{R}, f_c(\mathbf{x}) = F_c(\mathbf{x}, \mathbf{0}) \tag{2}$$

The idea is then to obtain ranking functions f_c from functions like F_c , as in (2), when F_c fulfils (1). Thus, given the set of preference judgments PJ_c , we can specify F_c by means of the constraints

$$\forall (\mathbf{v}, \mathbf{u}) \in PJ_c, F_c(\mathbf{v}, \mathbf{u}) > 0 \text{ and } F_c(\mathbf{u}, \mathbf{v}) < 0 \tag{3}$$

Therefore, PJ_c gives rise to a set of binary classification training set to induce F_c

$$E_c = \{(\mathbf{v}, \mathbf{u}, +1), (\mathbf{u}, \mathbf{v}, -1) : (\mathbf{v}, \mathbf{u}) \in PJ_c\} \tag{4}$$

Nevertheless, a separating function for E_c does not necessarily fulfill (1). Thus, we need an additional constraint. So, if we represent each object description \mathbf{x} in a higher dimensional feature space by means of $\phi(\mathbf{x})$, then we can represent pairs (\mathbf{x}, \mathbf{y}) by $\phi(\mathbf{x}) - \phi(\mathbf{y})$. Hence, a classification SVM can induce from E_c a function of the form:

$$F_c(\mathbf{x}, \mathbf{y}) = \sum_{s \in S(c)} \alpha_s z_s \langle \phi(\mathbf{x}_s^{(1)}) - \phi(\mathbf{x}_s^{(2)}), \phi(\mathbf{x}) - \phi(\mathbf{y}) \rangle \tag{5}$$

where $\langle \mathbf{x}, \mathbf{y} \rangle$ stands for the inner product of vectors \mathbf{x} and \mathbf{y} ; $S(c)$ is the set of support vectors, notice that they are formed by two d-dimensional vectors $(\mathbf{x}_s^{(1)}, \mathbf{x}_s^{(2)})$, while the scalars z_s represent the class +1 or -1. Trivially, F_c fulfils the condition (1). Let us remark that if k is a kernel function, defined as the inner product of two objects represented in the feature space, that is, $k(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$, then the kernel function used to induce F_c is

$$\mathbf{K}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) = k(\mathbf{x}_1, \mathbf{x}_3) - k(\mathbf{x}_1, \mathbf{x}_4) - k(\mathbf{x}_2, \mathbf{x}_3) + k(\mathbf{x}_2, \mathbf{x}_4) \tag{6}$$

Usually it is employed a linear or a simple polynomial kernel; that is, $k(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle$, or $k(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + 1)^g$, with $g = 2$.

Once we have a function F_c for a consumer c fulfilling (1), then, using (2), a ranking or preference or utility function f_c is given (but for an irrelevant constant) by

$$f_c(\mathbf{x}) = \sum_{s \in S(c)} \alpha_s z_s \langle \phi(\mathbf{x}_s^{(1)}) - \phi(\mathbf{x}_s^{(2)}), \phi(\mathbf{x}) \rangle = \sum_{s \in S(c)} \alpha_s z_s (k(\mathbf{x}_s^{(1)}, \mathbf{x}) - k(\mathbf{x}_s^{(2)}, \mathbf{x})) \tag{7}$$

Therefore, f_c can be represented by the weight vector \mathbf{w}^c in the higher dimensional space of features such that

$$f_c(\mathbf{x}) = \langle \mathbf{w}^c, \phi(\mathbf{x}) \rangle, \quad \mathbf{w}^c = \sum_{s \in S(c)} \alpha_s z_s (\phi(\mathbf{x}_s^{(1)}) - \phi(\mathbf{x}_s^{(2)})) \tag{8}$$

Notice that (8) defines the ranking of an object represented by a vector \mathbf{x} . This is not an absolute value; its importance is the relative position that gives to \mathbf{x} against to other objects \mathbf{y} in the *competition* for gaining the appreciation of consumer c . Now we only need to define the distance of consumers' preferences. Given that preferences are codified by those weighting vectors, we define the similarity of the preferences of consumer c and c' by the cosine of their weighting vectors. In symbols,

$$\text{similarity}(\mathbf{w}^c, \mathbf{w}^{c'}) = \cos(\mathbf{w}^c, \mathbf{w}^{c'}) = \frac{\langle \mathbf{w}^c, \mathbf{w}^{c'} \rangle}{\|\mathbf{w}^c\| * \|\mathbf{w}^{c'}\|} \tag{9}$$

Given that this definition uses scalar products instead of coordinates of weighting vectors, we can easily rewrite (10) in terms of the kernels used in the previous derivations. The essential equality is:

$$\begin{aligned} \langle \mathbf{w}^c, \mathbf{w}^{c'} \rangle &= \sum_{s \in S(c)} \sum_{s' \in S(c')} \alpha_s \alpha_{s'} z_s z_{s'} \langle \phi(\mathbf{x}_s^{(1)}) - \phi(\mathbf{x}_s^{(2)}), \phi(\mathbf{x}_{s'}^{(1)}) - \phi(\mathbf{x}_{s'}^{(2)}) \rangle \\ &= \sum_{s \in S(c)} \sum_{s' \in S(c')} \alpha_s \alpha_{s'} z_s z_{s'} \mathbf{K}(\mathbf{x}_s^{(1)}, \mathbf{x}_s^{(2)}, \mathbf{x}_{s'}^{(1)}, \mathbf{x}_{s'}^{(2)}) \end{aligned} \tag{10}$$

5 Clustering Consumers with Homogeneous Tastes

In the previous section we have associated one data point for each consumer in the space of preference criteria represented by ranking or preference functions. Moreover, we have defined a reasonable similarity measure for preference criteria; now we proceed to look for clusters of consumers with homogeneous tastes. For this purpose, we applied a nonparametric pairwise algorithm [8].

Let $S = (s_{ij})$ be a square matrix where s_{ij} stands for the similarity between data points i and j ; in our case, data points are the vectorial representation of the preference criteria of consumers, and similarities are given by equation (9). Then, matrix S is transformed iteratively, following a two step procedure that converges to a two values matrix (1 and 0), yielding a bipartition of the data set into two clusters. Then, recursively, the partition mechanism is applied to each of the resulting clusters represented by their corresponding submatrices. To guarantee that only meaningful splits take places, in [8] the authors provide a cross validation method that measures an index that can be read as a significance level; we will only accept splits which level is above 0.90.

The first step normalizes the columns of S using the L_∞ norm; then the proximities are re-estimated using the Jensen-Shannon divergence. The idea is to formalize that two preference criteria are close (after these two steps) if they were both similar and dissimilar to analogous sets of criteria before the transformation.

6 Experimental Results

In this section, we report the outputs obtained with the database of beef meat consumers. In order to consider significant opinions, we first selected those people involved in our consumers' panel whose ratings gave rise to at least 30 preference judgments; these yielded us to consider a set of 171 panelists that tested from 9 to 14 samples of meat of 101 different animals. The total amount of different samples was 303, since the meat from each animal was prepared with 3 different aging periods: 1, 7, and 21 days. Then the opinions of our panelists can be estimated inducing a preference or ranking function as was explained in section 4. Notice that only such functions can be used in order to compare the preferences of different consumers; in general, two arbi-

Table 2. For clusters of acceptance and tenderness datasets, this table reports the number of preference judgments (PJ), percentage of disagreements, and classification errors achieved into clusters with their own ranking or preference function, and using the function of the other cluster

Dataset	cluster	PJ	disagreements %	classification errors using function	
				own %	other %
acceptance	left	1927	16.19	19.20	50.96
	right	2150	17.07	21.12	54.95
tenderness	left	2487	15.96	19.38	61.98
	right	2432	15.21	19.59	61.06

trary consumers have not tested samples of the same animal prepared with the same aging. However, it is possible to compare the preference functions of any couple of consumers as vectors in a high dimension space following the kernel based method of section 4.

The clustering algorithm [8] returns the trees depicted in Figure 1. Split nodes achieved a confidence level of 91% for tenderness dataset, and 97% for acceptance. The leaves of these trees and the dataset of flavor reached lower confidence levels, and therefore they were rejected.

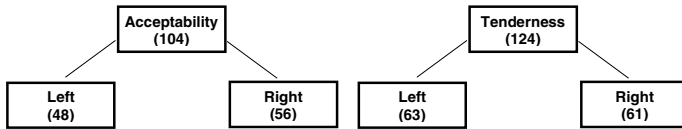


Fig. 1. Trace of the clustering algorithm. In each node we report the number of consumers

The job of clustering is to compute groups with minimal intra-group and maximal inter-group distances or differences. In our case, the relevance of clusters can be estimated, in part, by the coherence of consumers included into the same cluster, which can be measured by the classification error of the SVM used to compute the ranking or preference function of each cluster. Let us notice that the union of preference judgments of the members of the same cluster has some disagreements; if for each pair of samples we choose the most frequent relative ordering, then about 16% of preference pairs of each cluster express a particular disagreement with the majority opinion of the cluster, see Table 2. However, every preference judgment is included in the training set of each cluster; this sums more than 2000 preference judgments, what means (see equation 4 in section 4) more than 4000 training instances for the corresponding classification sets. When we use a polynomial kernel of degree 2, the errors range from 19.20% to 21.12%; we used this kernel following [2, 6, 7]. Nevertheless, if we apply the induced classification function of each cluster to the other one, then the errors rise to more than 50% in the case of acceptance, and more than 60% in the case of tenderness. Notice that in both cases we are ranking the same samples and these errors can be understood as the probability of reversing the order given by one of such clusters when we use the criteria of the other one. Therefore, 50% of error

means a random classification, and over that threshold means that ranking criteria is approaching the exactly opposite, see Table 2.

In general, it is well known that meat qualities are mainly the result of a set of complex factors. In this study, we are interested in knowing if there are different groups of people who prefer some breeds to others. To gain insight into the meaning of the preference criteria of each cluster, we used the ranking or preference functions to order the samples of meat; then we assessed 10 points to those samples included in the first decile, 9 to the second decile, and so on. Graphical representations of the average points obtained by each breed are shown in Figure 2; notice that the average score of all samples is 5.5. The results are quite the same if we use quartiles instead of deciles or any other division of the relative rankings of each cluster.

In the acceptance dataset (Fig. 2 left), let us emphasize the opposite role played by Retinta and Asturiana breeds: they were first and last (or almost last) in each cluster alternatively. In [6, 7] we used Boolean attributes to include the breed in the description of each sample, and then Retinta and Asturiana were found to be the most relevant Boolean features in order to explain consumer’s acceptance of meat. Additionally, these two breeds have significant differences in carcass composition (see Table 1). Notice that Asturiana breed is the only double muscled breed of the sample, and then it has the lowest values in percentages of subcutaneous and inter-muscular fat, and bone; while Retinta is the unimproved rustic breed with the highest percentages of fat and bone. Therefore, there are some reasons so as to assign opposite ratings to samples of these two breeds, although, in general, the final acceptance scorings rely on a complex set of features.

In tenderness dataset (Fig. 2 right), meat from Pirenaica and Retinta breeds are the tenderest for people in left cluster, however they are ranked in low positions in right cluster. We can say exactly the opposite of meat from Asturiana and Parda breeds. Again, Asturiana and Retinta breeds play opposites roles in each cluster.

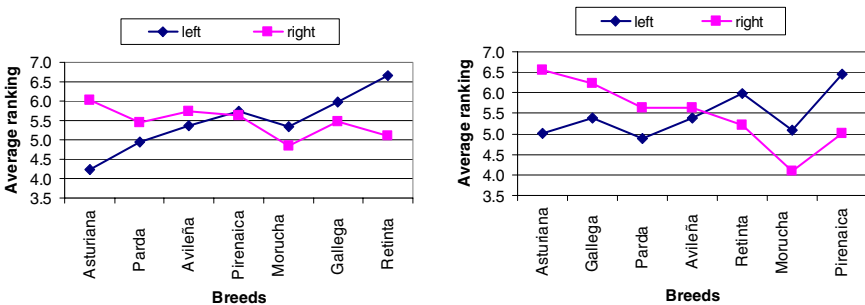


Fig. 2. Average ranking scores for each breed. Acceptance (left). Tenderness (right)

Acknowledgments

We would like to thank: the authors of Spider [9], a MatLab toolbox that includes kernel based algorithms; Thorsten Joachims for his SVM^{light} [10]. Those systems were used in the experiments reported in this paper; and INIA (Instituto Nacional de Inves-

tigación y Tecnología Agraria y Alimentaria of Spain) and Breeders Associations for the financial (Grant SC-97019) and technical support.

References

1. Sañudo, C.; Macie, E.S.; Olleta, J.L.; Villarroel, M.; Panea, B.; Albertí, P.. The effects of slaughter weight, breed type and ageing time on beef meat quality using two different texture devices. *Meat Science*, 66 (2004), 925–932
2. Díez, J.; Bayón, G. F.; Quevedo, J. R.; del Coz, J. J.; Luaces, O.; Alonso, J.; Bahamonde, A.. Discovering relevancies in very difficult regression problems: applications to sensory data analysis. *Proceedings of the European Conference on Artificial Intelligence (ECAI 2004)*, 993-994
3. Herbrich, R.; Graepel, T.; and Obermayer, K.: Large margin rank boundaries for ordinal regression. In A. Smola, P. Bartlett, B. Scholkopf, and D. Schuurmans, editors, *Advances in Large Margin Classifiers*, 115–132. MIT Press, Cambridge, MA, (2000)
4. Joachims, T.: Optimizing search engines using clickthrough data. In: *Proceedings of the ACM Conference on Knowledge Discovery and Data Mining (KDD) (2002)*
5. Bahamonde, A.; Bayón, G. F.; Díez, J.; Quevedo, J. R.; Luaces, O.; del Coz, J. J.; Alonso, J.; Goyache, F.. Feature subset selection for learning preferences: a case study. *Proceedings of the 21st International Conference on Machine Learning, (ICML 2004)*, 49-56
6. Del Coz, J. J.; Bayón, G. F.; Díez, J.; Luaces, O.; Bahamonde, A.; Sañudo, C.. Trait selection for assessing beef meat quality using non-linear SVM. *Proceedings of the Eighteenth Annual Conference on Neural Information Processing Systems (NIPS 2004)*, 321-328
7. Luaces, O.; Bayón, G.F.; Quevedo, J.R.; Díez, J.; del Coz, J.J.; Bahamonde, A.. Analyzing sensory data using non-linear preference learning with feature subset selection. *Proceedings of the 15th European Conference of Machine Learning, (ECML 2004)*, 286-297
8. Dubnov, S.; El-Yaniv, R.; Gdalyahu, Y.; Schneidman, E.; Tishby, N.; Yona, G.. A New Nonparametric Pairwise Clustering Algorithm Based on Iterative Estimation of Distance Profiles. *Machine Learning*, 47 (2002), 35–61
9. Weston, J.; Elisseeff, A.; BakIr, G.; Sinz, F.: SPIDER: object-orientated machine learning library. <http://www.kyb.tuebingen.mpg.de/bs/people/spider/>
10. Joachims, T.. *Making large-Scale SVM Learning Practical*. *Advances in Kernel Methods - Support Vector Learning*, B. Schölkopf and C. Burges and A. Smola (ed.), MIT-Press, (1999)

Corpus-Based Neural Network Method for Explaining Unknown Words by WordNet Senses

Bálint Gábor, Viktor Gyenes, and András Lőrincz*

Eötvös Loránd University, Pázmány P. sétány 1/C, Budapest 1117
{gbalint, gyenesvi}@inf.elte.hu, andras.lorincz@elte.hu
<http://nig.inf.elte.hu/>

Abstract. This paper introduces an unsupervised algorithm that collects senses contained in WordNet to explain words, whose meaning is unknown, but plenty of documents are available that contain the word in that unknown sense. Based on the widely accepted idea that the meaning of a word is characterized by its context, a neural network architecture was designed to reconstruct the meaning of the unknown word. The connections of the network were derived from word co-occurrences and word-sense statistics. The method was tested on 80 TOEFL synonym questions, from which 63 questions were answered correctly. This is comparable to other methods tested on the same questions, but using a larger corpus or richer lexical database. The approach was found robust against details of the architecture.

1 Introduction

The Internet is an immensely large database; large amount of domain specific text can be found. Intelligent tools are being developed to determine the meaning of documents, and manually created lexical databases are intended to provide help for such tools. However, manually assembled lexical databases are unable to cover specific, emerging subjects, thus documents may contain words of unknown meanings; words that are not contained in the lexical databases, or the contained meanings do not fit into the context found in the documents. However, the meaning of these words can often be inferred from their contexts of usage. The aim of our method is to explain words that are unknown to a human or machine reader, but are contained in many documents in the same sense.

To achieve this goal, our method looks for WordNet senses that are semantically close to the unknown meaning of the word. We rely on the common practice of measuring the similarity of words based on their contextual features, and designed a neural network architecture by means of three databases as sources of information. The first is WordNet¹, where the words are grouped into synonym sets, called *synsets*. The second source of information that our method exploits is SemCor², which is a corpus tagged with WordNet senses. SemCor was used

* Corresponding author

¹ <http://www.cogsci.princeton.edu/~wn/>

² <http://www.cs.unt.edu/~rada/downloads/semcor/semcor2.0.tar.gz>

to obtain information on the statistical distributions of the senses of the words. The third database used is the British National Corpus (BNC³), a collection of English texts of 100 million words. Our assumption is that the meaning of a word is similar to the meaning of a *sense*, if they appear in the same context, where we define context by the *senses* that they often co-occur with⁴.

Testing of any new method requires a controllable benchmark problem. Our method was evaluated on 80 synonym questions from the Test of English as a Foreign Language (TOEFL⁵). The system scored 78.75% (63 correct answers). Many other studies had also chosen this TOEFL benchmark problem: Landauer and Dumais's Latent Semantic Analysis (LSA) [1] is based on co-occurrences in a corpus, and it provides generalization capabilities. It was able to answer 64.4% of the questions correctly. Turney's Pointwise Mutual Information Information Retrieval (PMI-IR) algorithm [2] performed 73.25% on the same set of questions. This is also a co-occurrence based corpus method, which examines noun enumerations. It uses the whole web as a corpus and exploits AltaVista's special query operator, the *NEAR* operator. Terra and Clarke [3] compared several statistical co-occurrence based similarity measures on a one terabyte web corpus, and scored 81.25%. Jarmasz and Szapakovicz constructed a thesaurus-based method [4], which performed 78.75% on these questions. They utilized Roget's Thesaurus to calculate path lengths in the semantic relations graph between two words, from which a semantic similarity measure could be derived.

This paper is organized as follows: Section 2 details the neural network method. Section 3 describes the various test cases and presents the results. Discussion is provided in Section 4, conclusions are drawn in Section 5.

2 Methods

As it was already mentioned, our method exploits three databases (Fig. 1(A)). BNC is used to obtain word co-occurrence statistics. The following estimations are also required: given a synset, how frequently is one of its words used to express that sense, and, on the other hand, if a word is used, how frequently is it used in one of its senses. Database SemCor was used to obtain these statistics. Since all these pieces of information are needed, we only used words and synsets which occur in SemCor at least once. This means 23141 words and 22012 synsets. Later, we extended these sets by adding the trivial synsets, which contained single words. Then we could experiment with 54572 words and 53443 synsets.

2.1 Co-occurrence Measures

The aim of our method is to find semantically close synsets to an unknown word. Two words that occur in similar contexts can be considered as similar in

³ <http://www.natcorp.ox.ac.uk/>

⁴ The words *sense* and *synset* is used interchangeably in this paper.

⁵ <http://www.ets.org/>

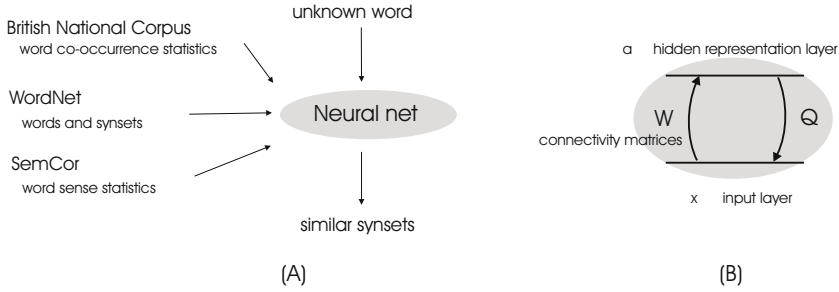


Fig. 1. Databases and reconstruction network architecture (A): Scheme of our method. (B): Basic computational architecture. Input layer (x) and hidden representation (a) are connected by bottom-up (W) and top-down (Q) matrices.

meaning. Therefore, we need to express the measures of co-occurrence between words and synsets; i.e. values indicating how often words and synsets co-occur.

The probability that word w_1 occurs near word w_2 can be estimated as follows: $P(w_1|w_2) = \frac{f(w_1, w_2)}{f(w_2)}$, where $f(w_1, w_2)$ is the number of times w_1 and w_2 co-occur in a 5 wide context window⁶ and $f(w)$ is the frequency of word w . We say that w_1 and w_2 are *near words* of each other, if both $P(w_1|w_2)$ and $P(w_2|w_1)$ are high, meaning that w_1 and w_2 are likely to co-occur. The following measure derived from mutual co-occurrences expresses this idea: $N(w_1, w_2) = \min(P(w_1|w_2), P(w_2|w_1))$. It is expected that this co-occurrence measure describes the contexts of the words. Given a word w , we call the *near word list* of w is the 100 words w_i for which the $N(w, w_i)$ values are the highest. This near word list can be represented as a *feature vector* of the word, the entries of the vector are the $N(w, w_i)$ values. Then the co-occurrence information about the words can be summarized in a quadratic and symmetric matrix N_W , where the i^{th} row of the matrix is the feature vector of the i^{th} word: $N_W(i, j) = N(w_i, w_j)$, where w_j is the j^{th} word in our vocabulary.

In SemCor, every occurrence of a word is tagged with a WordNet synset that expresses the meaning of the actual occurrence of the word. By counting these tags we can compute the desired probabilities. The probability that for a given word w , the expressed sense is s , can be estimated as $P(s|w) = \frac{f(w, s)}{f(w)}$, where $f(w, s)$ is the frequency of word w in sense s and $f(w)$ is the frequency of word w in any of its senses. We also need the probability that a given sense s is expressed by word w , which can be estimated as: $P(w|s) = \frac{f(w, s)}{f(s)}$, where $f(s)$ is the frequency of sense s , whichever word it is expressed by. These probabilities can also be summarized in matrix forms, denoted by S_W and W_S : $S_W(i, j) = P(s_i|w_j)$ and $W_S(i, j) = P(w_i|s_j)$.

Using the measures introduced, a co-occurrence measure between synsets can be derived. The idea is the following: given a synset s , the *near synsets* of s are the synsets of the near words of the words expressing s . This idea is

⁶ Increasing the context window by a factor of 2 had no significant effects.

expressed by the appropriate concatenation of the three matrices introduced above: $N_S = S_W N_W W_S$.

2.2 Reconstruction Networks

The basic reconstruction network model has two neuron layers. Connections bridge these layers. The lower layer is the input layer of the network, and the upper layer is called the hidden or internal representation layer (Fig. 1(B)). The network reconstructs its input by optimizing the hidden representation. For this reason, we call it reconstruction network. Formally, the following quadratic cost function is involved [5]:

$$J(a) = \|x - Qa\|_2^2, \quad (1)$$

where Q is the connectivity matrix, x is the input vector, a is the hidden representation vector. The columns of the connectivity matrix can be thought of as basis vectors, which must be linearly combined with the appropriate coefficients so that the combination falls close to the input. The optimization can either be solved directly

$$a = (Q^T Q)^{-1} Q^T x, \quad (2)$$

or iteratively

$$\Delta a \propto W(x - Qa), \quad (3)$$

where $W = Q^T$, which can be derived from the negative gradient of cost function (1). The form of (3) is more general than required by (1) but it still suitable as long as WQ is positive definite. Both methods have advantages and disadvantages. Directly solving the optimization returns the exact solution, but might require a considerable amount of memory, while the iterative solution requires less resources, but is computationally intensive.

The reconstruction network described above shall be called ‘one-tier’ network. We designed both ‘one-tier’ and ‘two-tier’ networks for the word-sense reconstruction. The two-tier network has two one-tiers on the top of each other. The internal representation layer of the first tier serves as input for the second tier. There are differences between the two architectures in computation speed and in numerical precision.

The feature vector representation of the context of the unknown word serves as input for the network. In the hidden layer, the neurons represent the candidate synsets. In the one-tier network (Fig. 2(A)), the top-down and bottom-up matrices are defined as follows: $Q = W_S N_S$ and $W = N_S^T S_W^T$. Thus, in the one-tier method, hidden synset activities are (a) transformed to near synset activities and then (b) the activities of the near words are generated. An illustrative iteration is depicted in Fig. 2(B): The activities in the topmost layer change during the reconstruction of the input word *frog*. It can be seen that only a few activities become high, others remain small. Note the horizontal scale: there were about 23,000 neurons in the topmost layer in this iteration.

⁷ However, we found that $\tilde{Q} = N_W W_S$ and $\tilde{W} = \tilde{Q}^T = W_S^T N_W^T$ are simpler, express the same relations and converge faster, thus these were used in the computations.

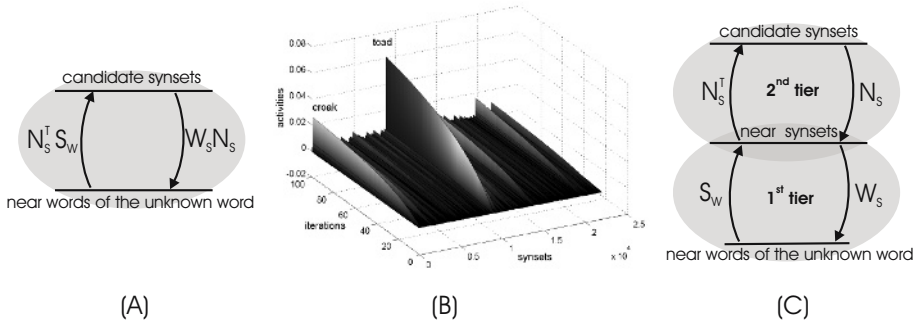


Fig. 2. One-tier and two-tier networks (A) and (C): one and two tier networks. The input of the i^{th} node of the lowest layer is $N(w, w_i)$, where w is the unknown word. $S_W(i, j) = P(s_i|w_j)$, $W_S(i, j) = P(w_i|s_j)$, $N_S = S_W N_W W_S$, where $N_W(i, j) = N(w_i, w_j)$ for all i, j . The *result* of the computation is the synset represented by the highest activity unit of the top layer after running the network. (B): Convergence of the iterative approximation. Input word is *frog*. High activity nodes correspond to sense *toad* and *croak*.

In the ‘two-tier’ network (Fig. 2(C)) the two steps of the transformation are separated. The nodes in the intermediate layer correspond to the near synsets of the unknown word. The bottom-up and top-down matrices are the S_W and the W_S matrices, respectively. In the second tier the connectivity matrix is N_S in both directions, which contains the synset near values. This captures the idea mentioned in the introduction; the meaning of the unknown word is similar to the meaning of a synset if they have the same near *synsets*.

3 Tests and Results

We tested our method on 80 TOEFL synonym questions, each question consisted of a *question word* (for example *grin*) and four *candidate answer words*, (for example: *exercise, rest, joke, smile*). The task was to find the candidate word that was the most similar in meaning to the question word (*smile*).

To meet our goals, we considered the question word the unknown word. We simulated the situation of the question word being *unknown* by erasing all information about the question word and its meaning from the SemCor statistics and WordNet synsets. After running the network for the context of the question word as input, we examined the activities corresponding to the synsets of the candidate words, and assigned a value to each candidate word equal to the highest activity of the synsets of the candidate word in the upper layer of the network. The candidate word with the highest value was the chosen answer.

In the one-tier network and in the second tier of the two-tier network, the huge number of connections between the nodes required the application of the less precise iterative method. However, the first tier of the two-tier network could be optimized directly, because the connectivity matrices were very sparse.

In some cases we have additional – top-down (TD) – information about the unknown word, for example its part-of-speech or the candidate answer words, this reduces the set of candidate synsets. The implementation of this filtering is simple in our system, since we can simply leave out the unnecessary synsets.

In order to test whether or not using synsets increases the efficiency of the system, we constructed a one tier *control* network which used only words. The input of the network was the same, however, the nodes in the upper layer corresponded to words. The connectivity matrix between the two layers is N_W , as defined in 2.1. The activities of the upper layer were examined; words having similar context as the input word were returned.

The number of correct answers in the various cases can be seen in Table 1. The best result, 63 correct answers (78.75%), was produced with part-of-speech constraint with the 23141 word data set and also without TD constraint with the 54572 word data set, utilizing the one-tier network. However, the best first iterations were achieved by the two-tier network. It can be seen, that iteration has improved the precision in almost all the cases. The control network started from 53 and 54 correct answers for the two word sets and reached 60 and 59 correct answers, respectively, by iterations. These results are considerably smaller than those of without the synsets, which supports our starting assumptions. We should note that if information related to the question words are not deleted from the database, then the number of correct answers is 68, which *amounts to 85%*.

Table 1. Results *No constr*: No constraint is applied. *PoS*: Synsets in the upper layer correspond to the part-of-speech of the word. *Candidate*: Only the synsets of the candidate words are used. *Control*: Single tier control network. *Direct*: non-iterative solution.

23141 words	No constr		PoS		Candidate		Control
	1 tier	2 tiers	1 tier	2 tiers	1 tier	2 tiers	1 tier
1 iter	55	58	55	58	55	58	53
10 iter	60	57	61	59	55	58	56
100 iter	61	60	63	58	58	56	60
direct	-	-	-	-	57	55	-
54572 words	No constr		PoS		Candidate		Control
	1 tier	2 tiers	1 tier	2 tiers	1 tier	2 tiers	1 tier
1 iter	56	59	56	59	56	59	54
10 iter	62	60	59	58	56	59	56
100 iter	63	61	62	58	57	56	59
direct	-	-	-	-	57	57	-

By examining the activities corresponding to the candidate answer words, decision points can be incorporated. Then the system may deny to answer a question, if the answer is uncertain. We could improve precision but the number of answers decreased considerably. Still, this property should be useful in multiple expert schemes, where experts may be responsible for different domains.

4 Discussion

Compared to the other methods, LSA performs relatively poorly (64.5%). However, the original intention of LSA was not to serve as an efficient TOEFL solver, but to model human memory. LSA reads the dictionary (the text database of the experiment) and runs the singular value decomposition only once and without knowing anything about the questions beforehand. After this procedure LSA can immediately answer the questions. While the first phase in LSA models a person's general learning process, this second phase imitates how someone solves questions without relying on any external aid [6]. By contrast, many other methods are allowed to use their databases after they have observed the questions.

Our method resembles LSA. Alike LSA, it works by the optimization of reconstruction using hidden variables over Euclidean norm. We also build a kind of memory model (the connectivity matrices of the neural network) before the questions are observed. When the questions are observed, the answer can be produced by running the network. Alike to LSA, our method was not developed for solving TOEFL questions, but for explaining unknown words. Considering this, our comparably high score (78.75%) is promising. True though, the original network incorporates information contained in WordNet and SemCor, however, the strength of the approach is shown by the *control network*, which did not use any lexical information, and gave 60 correct answers (i.e., 75%).

The Hyperspace Analogue to Language (HAL) model [7] works in high dimensions alike to our method. According to the HAL model, the strength of a term-term association is inversely proportional to the Euclidean distance between the context and the target words. Alike to HAL, our method makes use of the whole table of co-occurrences. This seems important; the larger table gave better result for us. Our method combines the advantages of LSA and HAL: it makes use of all information like HAL and adopts hidden variables like LSA.

We also included other information, the uncertainty of the answer, that goes beyond the statistics of co-occurrences. It may be worth noting here that our approach can be generalized to hidden, overcomplete, and sparse representations [8]. Such non-linear generalizations can go beyond simple computational advantages when additional *example based information* [9] or *supervisory training* are to be included.

A recent paper on meaning discovery using Google queries [10] thoroughly details the development of semantic distances between words. The method uses first order co-occurrence counts (Google page counts) to determine the semantic distance of two words. The article describes a semantic distance called Normalized Google Distance (NGD), derived from the same formula that we use to calculate the co-occurrence measure of two words. However, in our case, the formula was used to examine second order co-occurrences instead of first order co-occurrences. We conducted two studies with NGD. First, we solved the 80 TOEFL synonym questions using NGD as described in the paper; we measured the distance of the question word and each candidate word, and chosen the one with the smallest distance. Depending on the database we used to collect word frequencies, the results were different, however, surprisingly low: 30 correct an-

swers (37.5%) when BNC was used, and 40 correct answers (50.0%) when Google was used to return the page counts needed for NGD. In the other study we used our neural network method based on NGD instead of our co-occurrence measure. Results in this case were almost identical to the original setting, when we used our own co-occurrence measure, indicating the robustness of our solution against these details.

5 Conclusions

We have studied neural network architectures for explaining unknown words by known senses, senses that are contained in our lexical databases. We tested the method on TOEFL synonym questions. It was found that the networked solution provided good results, and was found robust against the details. At the cost of decreasing recall, the precision of the system can be improved. These features make our method attractive for various circumstances.

Acknowledgments

We are grateful to Prof. Landauer for providing us the TOEFL examples. We thank one of the referees for calling our attention to Burgess' HAL model.

References

- [1] Landauer, T.K., Dumais, S.T.: A solution to Plato's problem: The latent semantic analysis theory of acquisition, induction and representation of knowledge. *Psychol. Rev.* **104** (1997) 211–240
- [2] Turney, P.: Mining the Web for synonyms: PMI-IR versus LSA on TOEFL. In: *ECML Proceedings*, Freiburg, Germany (2001) 491–502
- [3] Terra, E., Clarke, C.L.A.: Choosing the word most typical in context using a lexical co-occurrence network. In: *Proc. of Conf. on Human Language Technol. and North American Chapter of Assoc. of Comput. Linguistics.* (2003) 244–251
- [4] Jarmasz, M., Szpakowicz, S.: Roget's Thesaurus and semantic similarity. In: *Proc. of the Int. Conf. on Recent Advances in Natural Language Proc. (RANLP-03).* (2003)
- [5] Haykin, S.: *Neural Networks: A comprehensive foundation.* Prentice Hall, New Jersey, USA (1999)
- [6] Landauer, T.K. (personal communication)
- [7] Burgess, C.: From simple associations to the building blocks of language: Modeling meaning in memory with the HAL model. *Behav. Res. Methods, Instr. and Comps.* **30** (1998) 188–198
- [8] Olshausen, B.A.: Learning linear, sparse factorial codes. *A.I. Memo 1580*, MIT AI Lab (1996) C.B.C.L. Paper No. 138.
- [9] Szatmáry, B., Szirtes, G., Lőrincz, A., Eggert, J., Körner, E.: Robust hierarchical image representation using non-negative matrix factorization with sparse code shrinkage preprocessing. *Pattern Anal. and Appl.* **6** (2003) 194–200
- [10] Cilibrasi, R., Vitanyi, P.: Automatic meaning discovery using google. *arXiv:cs.CL/0412098 v2* (2005)

Segment and Combine Approach for Non-parametric Time-Series Classification

Pierre Geurts, † and L. Wehenkel

University of Liège, Department of Electrical Engineering and Computer Science,
Sart-Tilman B28, B4000 Liège, Belgium,

† Postdoctoral researcher, FNRS, Belgium

{p.geurts,l.wehenkel}@ulg.ac.be

Abstract. This paper presents a novel, generic, scalable, autonomous, and flexible supervised learning algorithm for the classification of multivariate and variable length time series. The essential ingredients of the algorithm are randomization, segmentation of time-series, decision tree ensemble based learning of subseries classifiers, combination of subseries classification by voting, and cross-validation based temporal resolution adaptation. Experiments are carried out with this method on 10 synthetic and real-world datasets. They highlight the good behavior of the algorithm on a large diversity of problems. Our results are also highly competitive with existing approaches from the literature.

1 Learning to Classify Time-Series

Time-series classification is a challenging problem because of the high dimensionality and variable length of the data. Several classification methods have been proposed, including Gaussian discriminant analysis, support vector machines, and decision trees. However, these methods often require a large number of features and are sensitive to noise. In this paper, we propose a novel approach for time-series classification based on ensemble learning and temporal resolution adaptation.

First, the time-series are segmented into subseries of variable length. Then, a decision tree ensemble is trained on these subseries. The ensemble is trained using a bootstrap procedure, where each subseries is sampled with replacement. The final classification is obtained by voting over the ensemble members. This approach is robust to noise and can handle variable length time-series.

Our approach is evaluated on 10 datasets, including synthetic and real-world data. The results show that our method achieves high accuracy and is competitive with existing approaches. The proposed method is particularly effective for time-series classification tasks where the data is noisy and the length of the series varies significantly.

1 e , a 1 de , e ca 1 g) a d he . . . e hi di a ce , ea , e 1 c , bi a -
 1 . . . i h , ea , e , eighb he , e , e -ba ed , e h d [12,13]. A . . . e ia
 ad a , age , f he e a . . . ache 1 he . . . ibi 1 . . . bia he , e , e e a 1 , b e -
 . . . 11 g , 1 . . . be . . . eci c . . . edge. A he , a e 1 e , hi . . . be . . . eci c
 . . . de 1 g , e . . . a e he a . . . ica 1 . . . f , achi e ea , 1 g . . . a

The a . . . ach 1 , e 1 ga ed 1 . . . hi , a e , a 1 . . . a de e , 1 g a f . . . ge , e ic
 a d . . . - he , he f 1 e - e , e ca 1 ca 1 . . . e h d . M , e , eci e , he . . . ed
 a g , 1 h , e i e . . . a ge , e ic , e - . . . ce 1 g , age hich e , ac , f . . . he 1 e -
 . . . e i e a . . . be , f a d . . . e e c ed , b e i e , a , f he , a e e g h , hich
 a e abe ed 1 h he ca . . . f he 1 e - e , e f . . . hich he , e e a e . The a
 ge , e ic . . . e , 1 ed ea , 1 g , e h d 1 a . . . i ed . . . he , a . . . e f , b e i e , . . . a . . .
 de 1 ea , b e i e ca 1 e . F i a , a , e 1 e - e , e i ca 1 ed b , agge ga 1 g
 he , e dic 1 . . . fa 1 . . . b e i e , f he , a d 1 e . The e h d 1 c , b i ed 1 h
 a e - f d c . . . - a id a 1 . . . a e 1 . . . de . . . ad . . . a . . . a ica . . . he 1 e f
 he , b e i e . . . a gi e , da a e . A ba e ea , e , . . . e , e , e -ba ed , e h d
 beca e , f he , i ca ab i 1 a d a

Sec 1 . 2 , e e . . . a d . . . 1 a e he . . . ed a g , 1 h ic fa e . . . f , eg -
 . . . e a 1 . . . a d c , b i a 1 . . . f 1 e - e , e da a a d Sec 1 . 3 , e e . . . a e -
 . . . 1 ica e a a 1 . . . f he a g , 1 h . . . a di e , e e . . . f 1 e - e , e ca 1 ca 1 . . .
 a . . . F , he , de a 1 ab . . . hi . . . d . . . a be f , d 1 [4].

2 Segment and Combine

Notations. A 1 e - e , e 1 1 , i g i a . . . e , e e e d a a d i c , e e 1 e . . . 1 ed , a -
 1 . . . , ea - a ed , e c . . . , i g a . The di e e . . . c . . . e . . . f he e c . . . i g a a e
 ca ed e . . . a a , i b e 1 . . . ha f The . . . be , f 1 e - e , f , a gi e .
 e . . . , a a , i b e 1 ca ed 1 We e ha a e . . . , a a , i b e
 . . . f a gi e . 1 e - e , e ha e he a e d , a 1 . . . O , he , he ha d , he d , a 1 . . . f
 di e e . 1 e - e , e f a gi e . . . be (, da a e) a e . . . a . . . ed . . . be id e -
 ca . A gi e . 1 e - e , e 1 e a ed . . . a a , i c a , b e , a 1 . . . (. . . b e c) . A ea , -
 1 g , a . . . e (, da a e) 1 a e (, de 1 g 1 c . . . i de , ed 1 , e e a . . . a hi e e) f
 N , e c a 1 ed 1 e - e , e de . . . ed b $LS_N = \{ (\mathbf{a}(t^{d(o)}), o), c(o) \mid o = 1, \dots, N \}$,
 he , e o de . . . e a . . . b e , a 1 . . . , $d(o) \in \mathbb{N}$ a d f , he d , a 1 . . . f he 1 e -
 . . . e , e , $c(o)$, e fe . . . he ca . . . a . . . cia ed . . . he 1 e - e , e , a d

$$\mathbf{a}(t^{d(o)}, o) = (a_1(t^{d(o)}, o), \dots, a_n(t^{d(o)}, o))', a_i(t^{d(o)}, o) = (a_i(1, o), \dots, a_i(d(o), o))'$$

, e , e e . . . he e c . . . f n , ea - a ed e . . . a a , i b e 1 f d , a 1 . . . d(o).

The , b e c 1 e , f he 1 e - e , e ca 1 ca 1 b e 1 . . . de 1 e f . . . LS_N
 a ca 1 ca 1 . . . e c $\mathbf{a}(t^{d(o)}, o)$ hich , e dic ca e f a . . . e e 1 e -
 . . . e i e $\mathbf{a}(t^{d(o)}, o)$ a acc , a e a . . . i b e .

Training a subseries classifier. I 1 . . . , a 1 g , age , he , eg e . a d c , b i e
 a g , 1 h . . . e a 11 . . . a ba e ea , e . . . i ed a . . . b e i e ca 1 e f . . .
 LS_{N1} he f . . . 1 g , a :

Subseries sampling. For $i = 1, \dots, N_s$ choose $o_i \in \{1, \dots, N\}$, and, for each o_i , choose $t_i \in \{0, \dots, d(o_i) - \ell\}$, and, for each o_i , choose

$$\mathbf{a}_{t_i}^\ell(o_i) = (a_1(t_i + 1, o_i), \dots, a_1(t_i + \ell, o_i), \dots, a_n(t_i + 1, o_i), \dots, a_n(t_i + \ell, o_i))$$

and call the n -dimensional vector $\mathbf{a}_{t_i}^\ell(o_i)$ the i -th subseries of length ℓ . Collect the N_s subseries into the matrix

$$LS_{N_s}^\ell = \{(\mathbf{a}_{t_i}^\ell(o_i), c(o_i)) \mid i = 1, \dots, N_s\}.$$

Classifier training. Use the bagged ensemble $LS_{N_s}^\ell$ to build a bagged classifier. This classifier is denoted by $P_c^\ell(\mathbf{a}^\ell)$.

Notice that the N_s subseries have the same length ℓ , and the bagged classifier $P_c^\ell(\mathbf{a}^\ell)$ is trained on the matrix $LS_{N_s}^\ell$.

Classifying a time-series by votes on its subseries. For a test time-series $\mathbf{a}(t^{d(o)}, o)$, evaluate the classifier $P_c^\ell(\mathbf{a}^\ell)$ for each $o \in \{1, \dots, N\}$, and call the resulting

$$c(\mathbf{a}(t^{d(o)}, o)) \triangleq \arg \max_c \left\{ \sum_{i=0}^{d(o)-\ell} P_c^\ell(\mathbf{a}_i^\ell(o)) \right\}.$$

value $c(\mathbf{a}(t^{d(o)}, o))$ the o -th vote of the classifier $P_c^\ell(\mathbf{a}^\ell)$ on the test time-series $\mathbf{a}(t^{d(o)}, o)$.

Tuning the subseries length ℓ . In addition to the choice of bagging, the choice of ℓ is also important. In practice, the bagged classifier $P_c^\ell(\mathbf{a}^\ell)$ is trained on the matrix $LS_{N_s}^\ell$ and the test time-series $\mathbf{a}(t^{d(o)}, o)$ is classified by the classifier $P_c^\ell(\mathbf{a}^\ell)$. The choice of ℓ is also important because it affects the variance of the classifier. In practice, the bagged classifier $P_c^\ell(\mathbf{a}^\ell)$ is trained on the matrix $LS_{N_s}^\ell$ and the test time-series $\mathbf{a}(t^{d(o)}, o)$ is classified by the classifier $P_c^\ell(\mathbf{a}^\ell)$. The choice of ℓ is also important because it affects the variance of the classifier. In practice, the bagged classifier $P_c^\ell(\mathbf{a}^\ell)$ is trained on the matrix $LS_{N_s}^\ell$ and the test time-series $\mathbf{a}(t^{d(o)}, o)$ is classified by the classifier $P_c^\ell(\mathbf{a}^\ell)$.

Base learners. In this article, we consider the following base learners (SVM, k NN, MLP etc.) could be used to build the bagged classifier. However, for each base learner, we need to choose the parameters of the base learner. In this article, we consider the following base learners (SVM, k NN, MLP etc.) could be used to build the bagged classifier. However, for each base learner, we need to choose the parameters of the base learner.

Table 1. Summary of datasets

Dataset	Src.	N_d	n	c	$\underline{d} - \bar{d}$	Protocol	Best Ref	$\{\ell_i\}, \{s_i\}$
CBF	1	798	1	3	128	10-fold cv	0.00 [7]	1,2,4,8,16,32,64,96,128
CC	2	600	1	6	60	10-fold cv	0.83 [1]	1,2,5,10,20,30
CBF-tr	1	5000	1	3	128	10-fold cv	–	1,2,4,8,16,32,64,96,128
Two-pat	1	5000	1	3	128	10-fold cv	–	1,2,4,8,16,32,64,96,128
TTest	1	999	3	3	81-121	10-fold cv	0.50 [7]	3,5,10,20,40,60
Trace	3	1600	4	16	268-394	holdout 800	0.83 [1]	10,25,50,100,150,200,250
Auslan-s	2	200	8	10	32-101	10-fold cv	1.50 [1]	1,2,5,10,20,30
Auslan-b	5	2566	22	95	45-136	holdout 1000	2.10 [7]	1,2,5,10,20,30,40
JV	2	640	8	10	7-29	holdout 270	3.80 [8]	2,3,5,7
ECG	4	200	2	2	39-152	10-fold cv	–	1,2,5,10,20,30,39

¹<http://www.montefiore.ulg.ac.be/~geurts/thesis.html> ² [5] ³<http://www2.ife.no>
⁴<http://www-2.cs.cmu.edu/~bobski/pubs/tr01108.html>
⁵<http://waleed.web.cse.unsw.edu.au/new/phd.html>

de c, ibed 1, de a, i 1 [4]. I g, . . . a, ee b, ee c 1 g he be . . . i f, . . . a, a . . . e f ca, dida e, a d, . . . i. (b, h a, i b, e a d c - . . . i a, e, a d, i ed). Thi, e h d a, . . . ed ce, . . . g, a, i a, ce i h, i c, e a 1 g b i a, . . . ch. I i a, . . . i g i ca, . . . fa, e, i, he, a i g, age ha, baggi g, b, . . . i g, h i ch, . . . ea, ch f, . . . i a a, i b, e a d c - . . . i a, each, . . . de.

N, i ce ha, be ca, e he, eg, e, a d c, b i e a, . . . ach ha, . . . e i, i, i c, . . . a, i a, ce, ed c i, . . . ca a b i, . . . i i g, e a, . . . c, . . . e, . . . d c i e, . . . e, i g e, . . . ee i, h i c, . . . e. F, he a e, ea, . . . , he, . . . be, f, ee i, he, ee e, - e, be, e h d ca, be ch, e, . . . ea, . . . ab, . . . a (25 i, . . . e, e, i, e, . . .).

3 Empirical Analysis

3.1 Benchmark Problems

E, ee i, e, . . . a, e ca, ied, . . . 10, . . . be, . . . F, he a e, f, be i, . . . e, . . . ee, . . . i. Table 1 he, a i, . . . e, i e, f he 10 da a e, . We, efe, he i, e, e, ed, . . . eade, . . . [4] a d he, efe, e, ce, he, e i, f, . . . e de a i, . The, ec, d c, . . . g i e, . . . he (eb), . . . ce, f he da a e. The, e, f, . . . c, . . . g i e, he, . . . be, N_d , f i, e, e, i e, i, he da a e, he, . . . be, f, e, . . . a a, i b, e n , f each i, e, e, i e, he, . . . be, f, ca, e c , a d he, a g, e, f, a e, f he d, a i, $d(o)$; he, . . . e e, h c, . . . ec i e, . . . c, . . . de i e e, . . . a e; he i g h a d, i, h c, . . . g i e, e, ec i e, he be, . . . b i h e d, e, . . . a e (i h i d, e, i c a, . . . c, . . . a a b e, . . . c, . . .) a d he c, . . . e, . . . d i g, e, f, e, e, ce; he a, c, . . . g i e, he, i a, a e, ed f, he a a e e, ℓ a d s. The, . . . i, . . . be, . . . a e a, i - c i a, . . . be, . . . ec i ca, de i g, ed f, he a i d a i, . . . f i, e, e, i e, ca, i ca, i, . . . e h d, . . . h i e, he a, f, . . . be, . . . c, . . . e, . . . d, . . . ea, . . . d, . . . be, . . .

3.2 Accuracy Results

Acc, . . . ac, . . . e, . . . each, . . . be, . . . a, e g a, he, ed i, Table 2. I, . . . de, . . . a, e, . . . he i, e, e, f he, eg, e, a d c, b i e a, . . . ach, . . . ec, . . . a, e i, i h a, i, . . . e

Table 2. Error rates (in %) and optimal values of s and ℓ

Dataset	Temporal normalization				Segment&Combine ($N_s = 10000$)			
	ST		ET		ST		ET	
	Err%	s^*	Err%	s^*	Err%	ℓ^*	Err%	ℓ^*
CBF	4.26	24.0 ± 8.0	0.38	27.2 ± 7.3	1.25	92.8 ± 9.6	0.75	96.0 ± 0.0
CC	3.33	21.0 ± 17.0	0.67	41.0 ± 14.5	0.50	35.0 ± 6.7	0.33	37.0 ± 4.6
CBF-tr	13.28	30.4 ± 13.3	2.51	30.4 ± 4.8	1.63	41.6 ± 14.7	1.88	57.6 ± 31.4
Two-pat	25.12	8.0 ± 0.0	14.37	36.8 ± 46.1	2.00	96.0 ± 0.0	0.37	96.0 ± 0.0
TTest	18.42	40.0 ± 0.0	13.61	40.0 ± 0.0	3.00	80.0 ± 0.0	0.80	80.0 ± 0.0
Trace	50.13	50	40.62	50	8.25	250	5.00	250
Auslan-s	19.00	5.5 ± 1.5	4.50	10.2 ± 4.0	5.00	17.0 ± 7.8	1.00	13.0 ± 4.6
Auslan-b	22.82	10	4.51	10	18.40	40.0 ± 0.0	5.16	40.0 ± 0.0
JV	16.49	2	4.59	2	8.11	3	4.05	3
ECG	25.00	18.5 ± 10.0	15.50	19.0 ± 9.4	25.50	29.8 ± 6.0	24.00	32.4 ± 8.5

... a1a1... ech1e [2,6], hich a1... a... a... f... iga1e-e1e1... a... ec... f... ed di... e1... a1... f... ca... e1ca... a... ib... e: he1e1e... a... f... each... bec1... d1... d1... se... a... e... gh... eg... e... a... d... he... a... e... age... a... e... f... a... e... a... a... ib... e... a... g... he... e... eg... e... a... ec... ed, ied1... g... a... e... ec... f... n... s... a... ib... e... hich... a... e... ed... a1... he... ba... e... ea... e. The... a... a... ache... a... e... c... b1... ed... i... h... 1... g... e... dec1... ee (ST) a... d... e... e... be... f... 25... E... a... T... ee (ET) a... ba... e... ea... e. The... be... e... i... each... i... high... gh... ed.

... F... he... eg... e... a... d... c... b1... e... e... h... d... e... a... d... e... ac... ed... 10,000... b... e... 1e... The... i... a... a... e... f... he... a... a... e... e... ℓ ... a... d... s... a... e... ea... ched... a... g... he... ca... dida... e... a... e... e... ed... i... he... a... c... of... Table 1. Whe... he... e... 1... g... c... i... h... d... , he... a... a... e... a... e... ad... ed... b... 10-f... d... c... -... a... ida... he... ea... 1... g... a... e... ; he... he... e... 1... g... c... i... 10-f... d... c... -... a... ida... , he... ad... e... f... he... e... a... a... e... e... i... ade... f... each... f... he... e... f... d... b... a... 1... e... a... 10-f... d... c... -... a... ida... . I... hi... a... e... ca... e... a... e... age... a... e... a... d... a... da... d... de... ia... 1... f... he... a... a... e... e... s^* ... a... d... ℓ^* ... e... he... (e... e... a...)... e... 1... g... f... d... a... e... ed.

... F... he... e... e... e... be... e... ha... Seg... e... a... d... C... b1... e... i... h... E... a... T... ee (ET) ied... he... be... e... i... each... i... f... e... be... . O... he... he... be... (CBF, CBF-... , A... a... -b)... i... acc... ac... i... c... e... he... be... e... . O... he... ECG... be... , he... e... b... a... ed... a... e... e... ha... di... a... i... g... i... h... e... ec... he... a... 1... a... 1... a... ach... O... he... he... ha... d... i... c... ea... ha... he... c... b1... a... 1... f... he... a... 1... a... 1... ech1e... i... h... 1... g... ee (ST) i... e... a... ica... (ch)... e... acc... a... e... ha... he... he... a... ia... .

... We... a... be... e... ha... , b... h... f... a... 1... a... 1... a... d... eg... e... a... d... c... b1... e... , he... E... a... T... ee... a... a... g... i... e... ig... 1... ca... be... e... e... ha... 1... g... ee...¹... O... he... he... ha... d... he... 1... e... e... e... 1... g... f... he... eg... e... a... d... c... b1... e... e... h... d... i... ge... f... 1... g... e... dec1... ee... ha... f... E... a... T... ee... . I... deed... e... a... e... f... he... f... e... a... e... ed... ced... 1... a... e... age... b... 65%... h... i... e... a... e... f... he... a... e... a... e...

¹ There is only one exception, namely CBF-tr where the ST method is slightly better than ET in the case of “segment and combine”.

... ed ced b 30%. Ac a , ih eg e ad c. bie , ige ee ad E a-Tee a e c. e each he 1 e . f acc ac . ee a . be , hie he a e . ih . a i a i . Thi ca be e ai ed b he 1 i ic a ia ce ed c i . e ec f he eg e ad c. bie eh d, hichid e . he 1 a i c ea e f he ea i g a . e i e ad he a e agi g e ad . e ha i i ga e he e c a ia ce ed c i . ech i e i e e e be eh d (ee [3] f a d i c . i . f b i a a d a ia ce f he eg e ad c. bie eh d).

From the a e f $\ell^* 1$ he a c . . . f Tab e 2, i i cea ha he i a $\ell^* 1$ a . be de e de . a a e e . I deed, ih e ec . he a e age d a i . f he i e e i e hi . i a a e a ge f . 17% (. JV) . 80% (. TTe). Thi high h . he ef . e . f he a . a ic i g b c . a i da i . f $\ell^* a$ e a he ca ac i f he eg e ad c. bie a . ach . ada i ef . a i a b e e . a e . i . .

Ac a i . . f he e . . f he a . . c . . . f Tab e 2 ih he eigh h c . . . f Tab e 1, h . ha he eg e ad c. bie eh d ih E a-Tee i ac a . i e c . e i i e ih he be . b i hed e . . I deed, . CBF, CC, TTe , A a - , a d JV, i . e . a e e c . e . he be . b i hed . e .² Si ce . T ace, a d a e e e e . A a -b, he e . e e e . g d, e a a i de e e i e . . ee if he e i . . . f i . . e e . O T ace e e e a b e (ih E a-Tee a d $N_s = 15000$) . ed ce he e . . a e f . 5.00% . 0.875% b . . e a i g he i e e i e i . a . ed . be f 268 i e . i . . The a e a . ach ih 40 i e . i . . dec ea ed a . he e . . a e . A a -b f . 5.16% . 3.94%.

3.3 Interpretability

Le . i . . a e he . i b i i . e . ac i e . e a b e i f . a i . f . he . b e i e c a i e . Ac a , he e c a i e . . i d e f , each i e . i a ec . . e i a i g he c a - . b a b i i e . f . b e i e ce e ed a h i . i . He ce, . b e i e ha c . e . d . a high . b a b i i . f a ce, a i . b e f c a e ca be c . i d e ed a . i c a a e . . f h i . b e f c a e .

Fig e 1 . h . f . e a e , i . he . . a , . e . . a a i b e f . h ee i . a ce f he T ace . . be . e e c i e . f c a e 1, 3, a d 5, a d i . he b . . a . he e . i . . f he . b a b i i e . f he e h ee c a e a . ed i c e d f . . b e i e (f e g h $\ell = 50$) a he . . e . . g e i e f . . ef . . . gh . . he i e a i . The Ca . 3 i g a (. . i d d e) di e . f . . he Ca . 1 i g a (. . ef) . . i . he cc . . e ce f a . a i . i d a a e . i . . e f . he a i b e (a . . d $t = 200$); . . he . he ha d, Ca . 1 a d 3 di e . f . . Ca . 5 (. . i gh) i . he cc . . e ce f a . ha . ea i . he . he a i b e (a . . d $t = 75$ a d $t = 100$. e . e c i e) . From he . b a b i i . . . e . ee ha , f . $t \leq 50$ he h ee c a e a e e a . i e , b a . he i e he e he ea a ea . ($t \in [60 - 70]$) he . b a b i i . f Ca . 5 dec ea e f . he .

² Note that on CBF, CC, TTest, and Auslan-s, our test protocols are not strictly identical to those published since we could not use the same ten folds. This may be sufficient to explain small differences with respect to results from the literature.

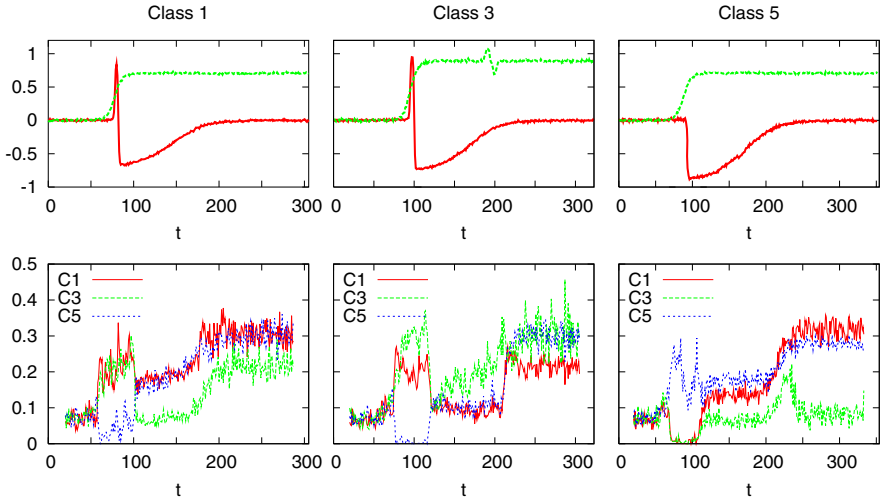


Fig. 1. Interpretability of “Segment and Combine” (Trace dataset, $N_s = 10000$, ET)

...erie (he.e.a.e.a.a.e.a.) a d i.c.e.a.e.f., he.igh-...erie (he.e...e.a.a.e.a.). S.b.e.e., a.d $t = 170$, he.b.e.r.e.1 he.i.d.d.e.1...a.c.e.a...d.e.c.he.1...i.d.a.a.e., hich.a.a.e.1...a.i.c.e.a.e.f.he...b.a.b.i.1...f.C.a.3, hie.f...he...he.1.e.e.r.e.C.a.e.1.a.d.5.b.e.c.e.e.a...1.e.a.d.C.a.3.e.a.r.e.e..N.ice.h.a.he...i.g.c.h.e.e...e.d...c.a.i.f...he.h.e.1.e.e.r.e.f...1...b.e.r.e.a...1.e.g.a.i.g...he.e.c...e.a.g...he.1.e.a.1.a.d.d.e.c.i.d.g...he...1.e.c.a...c.e.a...b.e.r.e.h.a.e.b.e.1.c...a.d.Thi...g.g.e...h...c.e.a...b.e.r.e.c.a.1.e.h.a.b.e.e...a.d, he.e.g.e...c...b.e.a...a.c.h.c.a.b.e...e.d...1...e.a-1.e.1...d.e...c.a.i.f...i.g.a...h...g...1.e.

4 Conclusion

I n h i s a r t i c l e, w e h a v e p r e s e n t e d a n e w g e n e r a l d a t a - a n a l y t i c a l m e t h o d f o r t i m e - s e r i e s c a t e g o r i e s w h i c h a d d i t i o n a l l y a c c o u n t f o r a g r e e m e n t s b e t w e e n f u t u r e t i m e - s e r i e s, i. d. e. a. l. b. e. r. e. c. a. l. e. f. o. r. h. i. s. a. r. t. i. c. l. e, a. d. c. a. l. e. t. i. m. e - s. e. r. i. e. b. a. e. a. g. g. h. e. m. e. d. i. c. i. n. e. t. i. m. e. b. e. r. e. T. h. e. b. e. r. e. e. g. h. i. s. a. n. a. l. y. t. i. c. a. l. a. d. a. d. e. d. b. y. t. h. e. a. g. r. i. h. t. h. e. e. n. t. i. t. l. e. 1. i. f. t. h. e. a. g. r. i. b. e. T. h. i. a. g. r. i. h. t. h. a. s. b. e. e. n. a. d. a. d. e. d. 10. b. e. c. h. a. n. g. e. b. e. r. e., h. e. r. e. i. n. e. d. e. d. e. n. t. i. c. a. l. e. i. n. e. i. t. h. a. e. f. f. h. e. a. a. g. r. i. h. t. f. o. r. t. h. e. i. e. a. r. e. G. i. v. e. t. h. e. d. i. e. r. e. n. t. f. o. r. b. e. c. h. a. n. g. e. b. e. r. e. a. d. c. e. s. a. n. i. n. c. i. d. e. n. t. f. o. r. a. g. r. i. h. t, h. i. s. a. r. t. i. c. l. e. i. n. t. r. o. d. u. c. e. s. F. o. r. t. h. e. r. e. s. t, t. h. e. d. i. s. t. r. i. b. u. t. i. o. n. e. a. c. c. o. u. n. t. i. n. g. t. i. m. e - s. e. r. i. e. h. a. s. b. e. e. n. h. i. g. h. l. i. g. h. t. e. d.

T. h. e. e. a. e. e. e. a. n. n. o. u. n. c. e. e. n. t. i. t. l. e. f. o. r. t. h. i. s. c. h. a. n. g. e. i. n. t. r. o. d. u. c. e. d. a. g. g. r. e. g. a. t. i. o. n. c. h. e. e. a. d. d. i. t. i. o. n. a. l. b. e. r. e. e. n. t. i. t. l. e. T. h. e. e. n. t. i. t. l. e. h. a. d. e. n. t. i. t. l. e. i. n. t. h. e. e. c. c. o. n. f. e. r. e. n. c. e. a. n. d. t. h. e. a. d. d. e. d. c. h. a. r. a. c. t. e. r. i. s. t. i. c. a. e. n. t. i. t. l. e.

of a label $l \in \mathcal{C}$. We have a \dots gged ha he e h d c d be ed
f ea - 1 e e e ie ca i ca i , b ad i g he i g che e.

The a ach e e ed he ef i e e ie i e e ia ide i ca he
e ed [9] f i age ca i ca i . Si i a idea c d a be e i ed
ied ge ic a ache f he ca i ca i f e b i g i ca e e ce .
A h gh he e a e be ha e di e e c a e ie , e be ie e
ha he e i b i f he a ach a e i i be ad i he e e
i a aigh f a d a e .

References

1. J. Alonso González and J. J. Rodríguez Díez. Boosting interval-based literals: Variable length and early classification. In M. Last, A. Kandel, and H. Bunke, editors, *Data mining in time series databases*. World Scientific, June 2004.
2. P. Geurts. Pattern extraction for time-series classification. In L. de Raedt and A. Siebes, editors, *Proceedings of PKDD 2001, 5th European Conference on Principles of Data Mining and Knowledge Discovery*, LNAI 2168, pages 115–127, Freiburg, September 2001. Springer-Verlag.
3. P. Geurts. *Contributions to decision tree induction: bias/variance tradeoff and time series classification*. PhD thesis, University of Liège, Belgium, May 2002.
4. P. Geurts and L. Wehenkel. Segment and combine approach for non-parametric time-series classification. Technical report, University of Liège, 2005.
5. S. Hettich and S. D. Bay. The UCI KDD archive, 1999. Irvine, CA: University of California, Department of Information and Computer Science. <http://kdd.ics.uci.edu>.
6. M. W. Kadous. Learning comprehensible descriptions of multivariate time series. In *Proceedings of the Sixteenth International Conference on Machine Learning, ICML'99*, pages 454–463, Bled, Slovenia, 1999.
7. M. W. Kadous and C. Sammut. Classification of multivariate time series and structured data using constructive induction. *Machine learning*, 58(1-2):179–216, February/March 2005.
8. M. Kudo, J. Toyama, and M. Shimbo. Multidimensional curve classification using passing-through regions. *Pattern Recognition Letters*, 20(11-13):1103–1111, 1999.
9. R. Marée, P. Geurts, J. Piater, and L. Wehenkel. Random subwindows for robust image classification. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR 2005)*, 2005.
10. I. Mierswa and K. Morik. Automatic feature extraction for classifying audio data. *Machine Learning*, 58(1-2):127–149, February/March 2005.
11. R. T. Olszewski. *Generalized feature extraction for structural pattern recognition in time-series data*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, 2001.
12. C.A. Ratanamahatana and E. Keogh. Making time-series classification more accurate using learned constraints. In *Proceedings of SIAM*, 2004.
13. H. Shimodaira, K.I. Noma, M. Nakai, and S. Sagayama. Dynamic time-alignment kernel in support vector machine. In *Advances in Neural Information Processing Systems 14, NIPS2001*, volume 2, pages 921–928, December 2001.
14. Y. Yamada, E. Suzuki, H. Yokoi, and K. Takabayashi. Decision-tree induction from time-series data based on standard-example split test. In *Proceedings of the 20th International Conference on Machine Learning (ICML-2003)*, 2003.

Producing Accurate Interpretable Clusters from High-Dimensional Data

Derek Greene and Padraig Cunningham

University of Dublin, Trinity College,
Dublin 2, Ireland

{derek.greene, padraig.cunningham}@cs.tcd.ie

Abstract. The primary goal of cluster analysis is to produce clusters that accurately reflect the natural groupings in the data. A second objective is to identify features that are descriptive of the clusters. In addition to these requirements, we often wish to allow objects to be associated with more than one cluster. In this paper we present a technique, based on the spectral co-clustering model, that is effective in meeting these objectives. Our evaluation on a range of text clustering problems shows that the proposed method yields accuracy superior to that afforded by existing techniques, while producing cluster descriptions that are amenable to human interpretation.

1 Introduction

The problem of generating meaningful clusters, a feature of the data, can be considered a challenging task. Figure 1, for example, identifies a set of clusters that accurately reflect the underlying structure of the data. A second objective is to identify features that are descriptive of the clusters. In addition to these requirements, we often wish to allow objects to be associated with more than one cluster. In this paper we present a technique, based on the spectral co-clustering model, that is effective in meeting these objectives. Our evaluation on a range of text clustering problems shows that the proposed method yields accuracy superior to that afforded by existing techniques, while producing cluster descriptions that are amenable to human interpretation.

The primary goal of cluster analysis is to produce clusters that accurately reflect the natural groupings in the data. A second objective is to identify features that are descriptive of the clusters. In addition to these requirements, we often wish to allow objects to be associated with more than one cluster. In this paper we present a technique, based on the spectral co-clustering model, that is effective in meeting these objectives. Our evaluation on a range of text clustering problems shows that the proposed method yields accuracy superior to that afforded by existing techniques, while producing cluster descriptions that are amenable to human interpretation.

The primary goal of cluster analysis is to produce clusters that accurately reflect the natural groupings in the data. A second objective is to identify features that are descriptive of the clusters. In addition to these requirements, we often wish to allow objects to be associated with more than one cluster. In this paper we present a technique, based on the spectral co-clustering model, that is effective in meeting these objectives. Our evaluation on a range of text clustering problems shows that the proposed method yields accuracy superior to that afforded by existing techniques, while producing cluster descriptions that are amenable to human interpretation.

g a 1 . . . ge e a e eigh . . . ig if 1 g he . . . e a ce f he e . . . 1 he c . . . cab a . . . each c . . . e, f . . . hich a e f c . . . e abe ca be de 1 ed. The . . . 11 . . . fd c . . . e eigh ca a . . . he hee d . . . e . . . ga a 1 igh 1 . . . a c . . . e 1 g . . . 1 . . . I a ic a , he a d c . . . e 1 a 1 g ed . . . a c . . . e , . . . e a 1 h . . . a if he c . . . de ce f he a 1 g e . . . Add 1 1 a , he e . . . f f c . . . e a e e e ca e he e a g i e d c . . . e . . . e a e . . . e ha . . . e . . . ic.

I h i a e , e 1 . . . d ce a c - c . . . e 1 g ech 1 e , ba ed . . . ec a a a - . . . 1 , ha . . . ide 1 e . . . e abe e be h i eigh . . . f b h e . . . a d d c - . . . e . . . F he . . . e , e h . . . ha b a . . . 1 ga 1 e a 1 e a 1 fac . . . 1 a 1 . . . che e , e ca . . . d ce a e e d c . . . e 1 g ha a . . . d 1 . . . ed acc . . . ac a d 1 e . . . e ab 1 1 . . . We c . . . a e . . . a g i h . . . 1 h e 1 1 g e h d . . . a a g e f da a e . . . a d b i e d i c . . . he ge e a 1 . . . f . . . e f c . . . e de c 1 1 . . . N e . . . ha a e e d e d e 1 . . . f h i a e 1 a a i a b e a a ech 1 c a . . . e . . . 1 h he a e 1 e [1].

2 Matrix Decomposition Methods

I h i e c 1 , e . . . e e a b i e f . . . a a f . . . e 1 1 g d i e 1 . . . ed c 1 . . . e h d ha ha e b e . . . e 1 . . . a 1 ed d c . . . e c . . . e 1 g . T de c i b e he a g i h . . . d i c e d 1 . . . he e a i d e f he a e , e e \mathbf{A} d e . . . e he $m \times n$ e . . . - d c . . . e . . . a 1 fac f n d c . . . e . . . , each f h i c h 1 e . . . e e d b a m -d i e . . . 1 a fea . . . e e c We a . . . e ha k 1 a 1 . . . a a e e . . . 1 d i c a 1 g he d e 1 ed . . . be . . . f c . . . e . . .

2.1 Spectral Co-clustering

S e c a c . . . e 1 g e h d ha e b e . . . ide . . . h ide a e e c 1 e . . . e a . . . f . . . d c i g d i 1 . . . a 1 1 . . . ac . . . a a g e f d a 1 . . . [2,3]. I . . . 1 . . . e . . . e . . . , he e a g i h . . . a a e he . . . ec a de c 1 1 . . . f a . . . a 1 . . . e . . . e . . . e 1 g a d a e 1 . . . de . . . c . . . e 1 . . . de 1 g . . . c . . . e . The ed ced . . . ace , c c e d f he e a d i g e i g e . . . e c 1 g a . . . e c f he . . . a 1 , ca b e i e d a a e . . . f e a i c a i a b e . . . a 1 g . . . 1 1 e . . . e g a 1 e a e .

A e a . . . each f . . . 1 . . . a e . . . c . . . e 1 g d c . . . e . . . a d e . . . a . . . g g e d b \mathbf{D} h i . . . [4] , he e he c - c . . . e 1 g . . . be . . . a f . . . a e d a he a . . . 1 a 1 . . . f he . . . 1 a a 1 ed c . . . o f a e i g h e d b i a 1 e g a h . I a . . . h . . . ha a a e a e d . . . 1 . . . a b e b a 1 e d b c 1 g he (SVD) of he a 1 $\mathbf{A}_n = \mathbf{D}_1^{-1/2} \mathbf{A} \mathbf{D}_2^{-1/2}$, he e $[D_1]_{ii} = \sum_{j=1}^n A_{ij}$ a d $[D_2]_{jj} = \sum_{i=1}^m A_{ij}$ a e diag . . . a . . . a i c e . A e d ced . . . e . . . e a 1 . \mathbf{Z}_1 he c c e d f he he f a d i g h 1 g a . . . e c f \mathbf{A}_n , c . . . e . . . d i g . . . he $g_2 k$ a g e 1 1 a 1 i g a . . . a e . V i e 1 g he . . . a 1 \mathbf{Z} a a l -d i e . . . 1 a g e . . . e i c e b e d d i g f he . . . i g 1 a d a a , he k - e a . . . a g i h . . . 1 a 1 e d 1 . . . h i . . . a c e d ce a d i 1 . . . c - c . . . e 1 g .

where P_i is a binary membership function. We define the k centroid of the cluster $\{\mu_1, \dots, \mu_k\}$.

After each iteration, we update the centroid of the data. According to [4], we have the centroid of the cluster μ_i is defined as follows. Let \mathbf{Z} and \mathbf{P} be the data matrix and the membership matrix:

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_1 \\ \mathbf{P}_2 \end{bmatrix}$$

where \mathbf{P}_1 and \mathbf{P}_2 denote the membership functions of each cluster. At each iteration, we update the centroid of the cluster μ_i as follows: $\mathbf{A}^T \hat{\mathbf{P}}_1$, where $\hat{\mathbf{P}}_1$ denotes the matrix \mathbf{P}_1 in the current iteration. This update is done by the centroid of the data \mathbf{Z} and the membership function. Similarly, we update the centroid of the cluster μ_j as follows: $\mathbf{V}^T \hat{\mathbf{P}}_2$, where $\hat{\mathbf{P}}_2$ denotes the membership matrix of the cluster μ_j in the current iteration.

Here, we have each iteration, we update the membership matrix \mathbf{P} by the following procedure. The membership function of the cluster μ_i is defined as follows: $\mu_i = \frac{S_{ij}}{\sum_l S_{il}}$, where S_{ij} is the membership function of the cluster μ_i in the current iteration. We can see that the membership function μ_i is in the range $[0, 1]$. We can see that the membership function μ_i is in the range $[0, 1]$. We can see that the membership function μ_i is in the range $[0, 1]$.

$$S_{ij} = \frac{1 + \cos(z_i, \mu_j)}{2}, \quad S_{ij} \leftarrow \frac{S_{ij}}{\sum_l S_{il}} \tag{1}$$

After the above procedure, we update the matrix \mathbf{S}_1 and \mathbf{S}_2 as follows: \mathbf{S}_1 is the matrix \mathbf{S}_1 in the current iteration and \mathbf{S}_2 is the matrix \mathbf{S}_2 in the current iteration.

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_1 \\ \mathbf{S}_2 \end{bmatrix}$$

Based on the above procedure, we update the membership matrix \mathbf{P} as follows: $\mathbf{A}^T \mathbf{S}_1$ and $\mathbf{A} \mathbf{S}_2$, where \mathbf{A} is the data matrix and \mathbf{S}_1 and \mathbf{S}_2 are the membership matrices of the clusters μ_i and μ_j in the current iteration.

3.2 Soft Spectral Co-clustering (SSC) Algorithm

Meanwhile, we define the co-clustering algorithm, where we use the co-clustering algorithm to find the clusters. The co-clustering algorithm is defined as follows: Let \mathbf{U} and \mathbf{V} be the data matrix and the membership matrix of the clusters μ_i and μ_j in the current iteration. We can see that the membership function μ_i is in the range $[0, 1]$. We can see that the membership function μ_i is in the range $[0, 1]$.

has the following nice properties. It is a low-rank approximation of the ground truth matrix \mathbf{P} . We believe that the proposed method will be useful for a variety of applications.

The proposed method is based on the following idea. We consider the matrix \mathbf{P} as a sum of two matrices. The first matrix is a low-rank approximation of \mathbf{P} and the second matrix is a matrix that captures the remaining information. This is done by solving the following optimization problem:

$$\min_{\mathbf{U}, \mathbf{V}} \|\mathbf{P} - \mathbf{U}\mathbf{V}^T\|_F$$

where \mathbf{U} and \mathbf{V} are matrices of size $n \times k$ and $k \times n$ respectively. This problem can be solved using the singular value decomposition (SVD) of \mathbf{P} . Let $\mathbf{P} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be the SVD of \mathbf{P} , where \mathbf{U} and \mathbf{V} are orthogonal matrices and $\mathbf{\Sigma}$ is a diagonal matrix. Then the optimal solution to the above problem is given by $\mathbf{U}_k = \mathbf{U}\mathbf{\Sigma}_k^{-1/2}$ and $\mathbf{V}_k = \mathbf{V}\mathbf{\Sigma}_k^{-1/2}$, where $\mathbf{\Sigma}_k$ is the k -th largest singular value of \mathbf{P} .

1. Compute the k -th largest singular value of \mathbf{A}_n and the corresponding left and right singular vectors $\mathbf{U}_k = (u_1, \dots, u_k)$ and $\mathbf{V}_k = (v_1, \dots, v_k)$.
2. Compute the low-rank approximation $\mathbf{Z} = \mathbf{U}_k \mathbf{V}_k^T$.

$$\mathbf{Z} = \begin{bmatrix} \mathbf{D}_1^{-1/2} \mathbf{U}_k \\ \mathbf{D}_2^{-1/2} \mathbf{V}_k \end{bmatrix}$$

3. Approximate the k -th largest singular value of \mathbf{Z} and the corresponding left and right singular vectors \mathbf{S}_1 and $\hat{\mathbf{P}}_1$.
4. Compute the final low-rank approximation $\mathbf{U} = \mathbf{A}\hat{\mathbf{P}}_2$ and $\mathbf{V} = \mathbf{A}^T \mathbf{S}_1$.

We note that the proposed method is a generalization of the method proposed in [2], where each cluster is assumed to be a cluster of size 90° from the set of clusters. However, the proposed method is able to handle clusters of arbitrary size and shape. The proposed method will be useful for a variety of applications, including image and video clustering.

4 Refined Soft Spectral Clustering

We now describe the refined soft spectral clustering algorithm. The algorithm is based on the following idea. We consider the matrix \mathbf{P} as a sum of two matrices. The first matrix is a low-rank approximation of \mathbf{P} and the second matrix is a matrix that captures the remaining information. This is done by solving the following optimization problem:

$$\min_{\mathbf{U}, \mathbf{V}} \|\mathbf{P} - \mathbf{U}\mathbf{V}^T\|_F$$

where \mathbf{U} and \mathbf{V} are matrices of size $n \times k$ and $k \times n$ respectively. This problem can be solved using the singular value decomposition (SVD) of \mathbf{P} . Let $\mathbf{P} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be the SVD of \mathbf{P} , where \mathbf{U} and \mathbf{V} are orthogonal matrices and $\mathbf{\Sigma}$ is a diagonal matrix. Then the optimal solution to the above problem is given by $\mathbf{U}_k = \mathbf{U}\mathbf{\Sigma}_k^{-1/2}$ and $\mathbf{V}_k = \mathbf{V}\mathbf{\Sigma}_k^{-1/2}$, where $\mathbf{\Sigma}_k$ is the k -th largest singular value of \mathbf{P} .

The refined soft spectral clustering algorithm is based on the following idea. We consider the matrix \mathbf{P} as a sum of two matrices. The first matrix is a low-rank approximation of \mathbf{P} and the second matrix is a matrix that captures the remaining information. This is done by solving the following optimization problem:

$$\min_{\mathbf{U}, \mathbf{V}} \|\mathbf{P} - \mathbf{U}\mathbf{V}^T\|_F$$

where \mathbf{U} and \mathbf{V} are matrices of size $n \times k$ and $k \times n$ respectively. This problem can be solved using the singular value decomposition (SVD) of \mathbf{P} . Let $\mathbf{P} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be the SVD of \mathbf{P} , where \mathbf{U} and \mathbf{V} are orthogonal matrices and $\mathbf{\Sigma}$ is a diagonal matrix. Then the optimal solution to the above problem is given by $\mathbf{U}_k = \mathbf{U}\mathbf{\Sigma}_k^{-1/2}$ and $\mathbf{V}_k = \mathbf{V}\mathbf{\Sigma}_k^{-1/2}$, where $\mathbf{\Sigma}_k$ is the k -th largest singular value of \mathbf{P} .

1. The matrix A is decomposed into a sum of k rank-1 matrices. Each rank-1 matrix is then normalized to have a unit norm. The resulting matrix is then used to initialize the clustering algorithm.

4.1 Refined Soft Spectral Co-clustering (RSSC) Algorithm

In the SSC algorithm described in Section 3.1, the choice of the number of clusters is a critical parameter. Here, we propose a refined algorithm, RSSC, which is based on the idea of soft clustering. The algorithm is as follows:

We define the matrix U and V by the following equations:

$$D(A||UV^T) = \sum_{i=1}^m \sum_{j=1}^n \left(A_{ij} \log \frac{A_{ij}}{[UV^T]_{ij}} - A_{ij} + [UV^T]_{ij} \right) \quad (2)$$

This function can be minimized using the Kullback-Leibler divergence. The algorithm is as follows:

1. Decompose A_n and cluster the embedded space Z as described in Section 3.1.
2. Assign k clusters to the rows of Z and calculate the initial matrices S_1 and S_2 as described.
3. Generate the initial factors $U = A^T S_1$ and $V = A S_2$.
4. Update V using the following equation:

$$v_{ij} \leftarrow v_{ij} \left[\left(\frac{A_{ij}}{[UV^T]_{ij}} \right)^T U \right]_{ij} \quad (3)$$

5. Update U using the following equation:

$$u_{ij} \leftarrow u_{ij} \left[\frac{A_{ij}}{[UV^T]_{ij}} V \right]_{ij}, \quad u_{ij} \leftarrow u_{ij} \frac{U_{ij}}{\sum_{l=1}^m U_{lj}} \quad (4)$$

6. Repeat steps 4 and 5 until convergence.

The algorithm is able to handle high-dimensional data by using the idea of soft clustering. The algorithm is as follows:

5 Experimental Evaluation

In this paper, we evaluate the accuracy of the SSC and RSSC algorithms using the standard cosine similarity (CC) based clustering algorithm and NMF algorithm. The datasets used are the 20 newsgroups (2) and the 11 newsgroups. The clustering algorithm used is the k means algorithm, which is based on the cosine similarity. For the clustering algorithm, we use the k means algorithm to cluster the data.

The evaluation is based on the accuracy of the clustering algorithm, which is the ratio of the number of correctly classified documents. For the clustering algorithm, we use the cosine similarity [1]. The evaluation is based on the accuracy of the clustering algorithm. We believe the clustering algorithm has high accuracy. NMF, the feature extraction algorithm, is also evaluated.

5.1 Results

The clustering accuracy, evaluated using the Normalized Mutual Information (NMI) evaluation algorithm [9]. We used the standard cosine similarity based clustering algorithm to cluster the data. The clustering algorithm used is the k means algorithm. The evaluation is based on the accuracy of the clustering algorithm. We believe the clustering algorithm has high accuracy. NMF, the feature extraction algorithm, is also evaluated.

Table 1 shows the performance of the clustering algorithm on the 20 newsgroups dataset. In general, the accuracy of the clustering algorithm is high. The clustering algorithm used is the k means algorithm. The evaluation is based on the accuracy of the clustering algorithm. We believe the clustering algorithm has high accuracy. NMF, the feature extraction algorithm, is also evaluated.

Table 1. Performance comparison based on NMI

Dataset	CC	NMF	SSC	RSSC
bbc	0.78	0.80	0.82	0.86
bbcsport	0.64	0.69	0.65	0.70
classic2	0.29	0.34	0.46	0.79
classic3	0.92	0.93	0.92	0.93
classic	0.63	0.70	0.62	0.87
ng17-19	0.39	0.36	0.45	0.50

Dataset	CC	NMF	SSC	RSSC
ng3	0.68	0.78	0.70	0.84
re0	0.33	0.39	0.35	0.40
re1	0.39	0.42	0.41	0.43
reviews	0.34	0.53	0.40	0.57
tr31	0.38	0.54	0.51	0.65
tr41	0.58	0.60	0.67	0.67

... d ce, he ea he de e. i i ic a e f he i i a i a i ... a eg e ... ed
 b he e ... ed a g i h ... ead ... ab e ... i ...

5.2 Cluster Labels

Gr e he e ... e be h i ... eigh ... d ced b he SSC a d RSSC a g i h ...
 a a a a a ... ach ... ge e a i g a e ... f ab e ... f each c ... e ... e ec
 he e ... i h he high ... a e f ... each c ... f he a ... U. D e ...
 ... ace e ... ic i ... , e ide a a ... e f he ab e ... e ec ed f ... c ... e ...
 ... d ced b he RSSC a g i h ... he ... da a e i Table 2, he e he a ... a
 ca eg ... a e: b ... i e ... , ... , e ... a ... e ... a d ech ... g .

Table 2. Labels produced by RSSC algorithm for *bbc* dataset

Cluster	Top 7 Terms
C1	company, market, firm, bank, sales, prices, economy
C2	government, labour, party, election, election, people, minister
C3	game, play, win, players, england, club, match
C4	film, best, awards, music, star, show, actor
C5	people, technology, mobile, phone, game, service, users

6 Concluding Remarks

I n h i a e , e de c i b e d a ... e h d b a e d ... e c , a a a ... i h a c a ... i e d
 ... a b e i e ... e a b e c ... e ... i ... a ... e high-d i ... e i ... a ... a c e . S b e e ... , e
 ... i ... d ced a ... e a ... ach ... a c h e e a ... e acc ... a e c ... e i g b a ... i g
 a c ... a i e d ... a ... i f a c ... i a ... i a ... c h e e ... e ... e a ... i i a ... i ... d ced
 ... i g ... e c ... a ... e c h ... i e . E a a ... i ... c ... d ced ... a ... a ... i e ... f e ... c ... a
 d e ... a e h a h i ... e h d c a ... e a d ... he i ... d ... e d i d e ... i c a ... i ... f ... e -
 a ... i g c ... e , h i e ... i ... a e ... d c i g d c ... e ... a d e ... e i g h ...
 h a a e a e a b e ... h ... a ... i e ... e a ... i ...

References

- Greene, D., Cunningham, P.: Producing accurate interpretable clusters from high-dimensional data. Technical Report CS-2005-42, Trinity College Dublin (2005)
- Ng, A., Jordan, M., Weiss, Y.: On spectral clustering: Analysis and an algorithm. In: Proc. Advances in Neural Information Processing. (2001)
- Brand, M., Huang, K.: A unifying theorem for spectral embedding and clustering. In: Proc. 9th Int. Workshop on AI and Statistics. (2003)
- Dhillon, I.S.: Co-clustering documents and words using bipartite spectral graph partitioning. In: Knowledge Discovery and Data Mining. (2001) 269–274
- Lee, D.D., Seung, H.S.: Learning the parts of objects by non-negative matrix factorization. Nature **401** (1999) 788–91
- Xu, W., Liu, X., Gong, Y.: Document clustering based on non-negative matrix factorization. In: Proc. 26th Int. ACM SIGIR. (2003) 267–273

7. Dhillon, I.S., Modha, D.S.: Concept decompositions for large sparse text data using clustering. *Machine Learning* **42** (2001) 143–175
8. Zhao, Y., Karypis, G.: Soft clustering criterion functions for partitional document clustering: a summary of results. In: *Proc. 13th ACM Conf. on Information and Knowledge Management*. (2004) 246–247
9. Strehl, A., Ghosh, J.: Cluster ensembles - a knowledge reuse framework for combining multiple partitions. *JMLR* **3** (2002) 583–617

Stress-Testing Hoeffding Trees

Geoff H. E., Richard Kirkby, and Bernhard Pfahringer

Department of Computer Science,
University of Waikato,
Hamilton, New Zealand

{geoff, rkirkby, bernhard}@cs.waikato.ac.nz

Abstract. Hoeffding trees are state-of-the-art in classification for data streams. They perform prediction by choosing the majority class at each leaf. Their predictive accuracy can be increased by adding Naive Bayes models at the leaves of the trees. By stress-testing these two prediction methods using noise and more complex concepts and an order of magnitude more instances than in previous studies, we discover situations where the Naive Bayes method outperforms the standard Hoeffding tree initially but is eventually overtaken. The reason for this crossover is determined and a hybrid adaptive method is proposed that generally outperforms the two original prediction methods for both simple and complex concepts as well as under noise.

1 Introduction

The Hoeffding tree algorithm [2] has been well established as the leading method for data stream classification. Several Hoeffding tree variants have been proposed, each with different predictive performance. Pfeiffer et al. [3] have shown how adding a Naive Bayes model to the Hoeffding tree can improve its performance, and we have also shown that the Naive Bayes model can be used to improve the performance of the Hoeffding tree.

In this paper we evaluate the performance of the Hoeffding tree and the Naive Bayes model in a variety of situations. We show that the Naive Bayes model can improve the performance of the Hoeffding tree in some situations, but that it can also be outperformed by the Hoeffding tree in other situations. We also show that the Naive Bayes model can be used to improve the performance of the Hoeffding tree in some situations, but that it can also be outperformed by the Hoeffding tree in other situations. We also show that the Naive Bayes model can be used to improve the performance of the Hoeffding tree in some situations, but that it can also be outperformed by the Hoeffding tree in other situations.

The paper is organized as follows. Section 2 contains a description of the Hoeffding tree algorithm. Section 3 describes the Naive Bayes model. Section 4 evaluates the performance of the Hoeffding tree and the Naive Bayes model. Section 5 concludes the paper.

2 Examining Hoeffding Trees

Da a , ea , ee , i e , e , e f , e a a i , d e , h e , e f da a a a b e a d h e a - i e , e f h e a g i h , d e e a i a i . We c o i d e a e h d f e a a i , h a e , i , h i , e , h i , a i i i g e f h e d a a . Th i a c h i e d b i g e e , i a c e a a e i g e a e e h e c , e , d e b e f , e i g i , a i h e , d e , i c e e a d a i g a i i c a e a c h i .

The a i c a i , e e a i , f H e d i g T e e i d c i , d i c e d i , h i a e , e i f , a i , g a i a h e , i c i e i , h e i g i a V F D T H e d i g b i d f , a i , [2] d e e i e h e i (i g a a e e , $\delta = 10^{-6}$, $\tau = 5\%$, a d $n_{min} = 300$) , a d h a d e , e i c a , i b e b G a i a a , - i a i , (h , g h , h t , e f e , h i a g i h a d h t n b h e a e a g i h i h N a i e B a e , e d i c i , a h e e a e) .

O , e e , a a i , a h e d i e e c e i a c c , a c b e e e h t a d h t n b . We a i h d a g e e a e d b a a d , c e , c e d d e c i i , e e c i i g f 10 , i a a , i b e i h 5 a e e a c h , 10 , e i c a , i b e , 2 c a e , a , e e d e h f 5 , i h e a e , a i g a e e 3 a d a 0.15 c h a c e f e a e h e a f e (h e a , e e h a d 741 d e , 509 f h i c h e e e a e) h i c h e h a , e f e , a h e i e , a d , e e . N e h a f , a g a h i h i a e e h a e a e a g e d , e , 10 , e e i i a e , d e e e .

F i g . e 1 . h e , e f e a a i g , e , 10 , i i , i a c e i h e , e e . A i , e i , d i e , i i c e a , h a h t n b g i e a i , e e i c a i c a i , a c c , i h b , h a i a , e f , i g e , e a c h i g a , d 99% a c c , a c i h e , g , .

F i g . e 2 . h e i a c h a , i e h a . 10% i i e a i , d c e d , h e d a a i h i f , a d , e e . A d i e e i c , e e e g e h t n b , e e i i i a , b , e h e b e f , e 2 . i i , i a c e h e g a h c , e a d i h e , g , h t n b i , e e .

N e , h e , e e g e e a , i a d , e d , d c e a c , e , a d , e e 50 , i a a , i b e i h 5 a e e a c h , 50 , e i c a , i b e , 2 c a e , a , e e

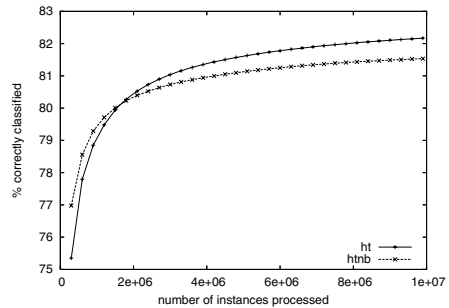
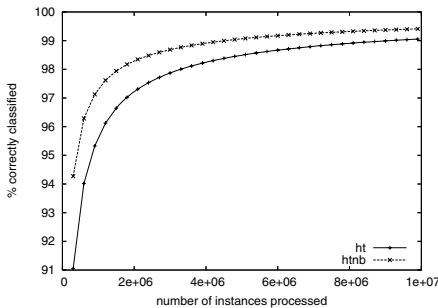


Fig.1. Simple random tree generator with no noise

Fig.2. Simple random tree generator with 10% noise

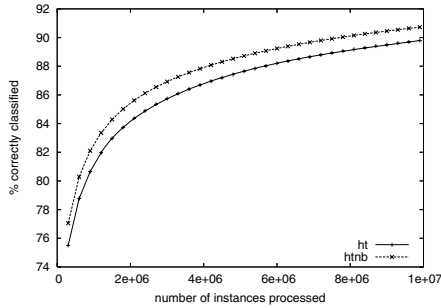


Fig. 3. Complex random tree generator with no noise

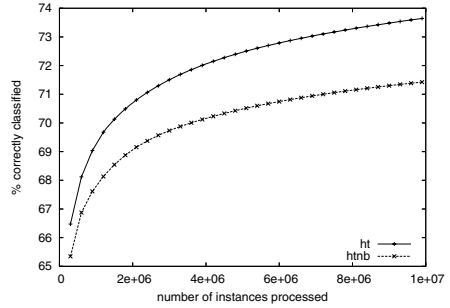


Fig. 4. Complex random tree generator with 10% noise

de h.f10, 1 h ea e. a 1 g a e e 5 a d a 0.15 cha ce f ea e h e eaf e. (he a e e e had 127,837 de , 90,259 f hich e e ea e).

Fig e 3 a d 4 h h e ea a i g c e e e i g f e h e c e e a d . e e da a , b h ce a a d e . The a e a e i b e e d , h h i e htnb i e e f e e h e e e .

Te i g h e i e e e i h h e e i e e e g e e e i i a e . Fig e 2 , a 5% i e h e c e e e e e e e e a e , a d a 20% i e c e e e a e h a h e 10% ca e . O h e c e e e e e da a h e g a i Fig e 4 i d e a h e e i e i ce e e .

3 Possible Solutions to the Problem

The e a a i e c e e d ca e h e e htnb i e e a c e a e ha ht . Thi be ha i , a ea e h e e i e i i d ced , a d be e e e e e e e e ced h e h e c e e i e e d i c e e ea e . The e be h e e e e h e da a e a e , h g h e e i e i e e a a e d e h e i i e f i a ce ha a e e e d be e c e e e e e e e e .

We e ca e ha h e e be i d e e e a d i c e i h e e e i c e bi a i i h i i da a . The ea e e e e a i i e e e e be e a e e , a d h e e h e d e e a e e i , ca i g htnb i e e e a c e a e ha ht .

E e e i e e i h e e e e i i e d a i e e e e e e e e e d h i h e e i : htnb e e e a e e e e f e e d ht i e ce e e g e h a a i c i a e e e d .

If h e e be i h htnb i d e e i e i d e e h e d e e i g e e i a b e i h e i ea e e e e , h e h e e a e e e a a h e e be c d be e e d .

O e e i h a be e e e e d b G a a h i [3] . g g e h e e e f a h e e e e e e e e e h i g e e f h e e e e e e i a ce . A e be i h h i i i i d e e i i g a e e e e e f h e e a h e e e g e e i c e e e e i , f e e f h e i a ce i h e cache i be a i ca b e e h e e ea e d e e h e e e (e e f e e h i a htnb-stm x h e e x i h e cache i e) .

A h e i d e a i i i h e i i f a a i f e e a e e e e e a i h a be e c h e e . F a a a i b e e i , i i i e i b e e a e e i a e h e d i b i i f

a e e i g f . he . i . F . he . he a i b e , e . i f . a i . i . . .
 ab . he e f he e . i , b e ca a . e ha he d i b i . i he
 a e a i he a e . Thi a . i a i . a beg . . . i c . . ec , b a e a
 i g i e he . de a a i g i i . a he ha . a i g i h . i f . a i .
(htnbp) .

A . e i a . e i . . . be i h h i a . ach i ha i f he a i ic ed
 . a e . i ch ice a e . i ed he . he . i deci . . . i be a e ed , ha i g
 a i ac . . he . ee . c . e . A . i be h . i Sec i . 4 , h i ge e a
 ha a de . i e . a e ec . . acc . ac .

A . . . i . . . h i . . . a i a i a e a a e . . de . e eaf ha i . ed f .
 . edic i . . . e . . , a d ea e he . i deci . . . a i ic . . . ched . Thi
 e ec i e d be he . . age e i e e . e eaf **(htnbps)** .

A . ada i e . . . i . . . ee h . f e . he Nai e Ba e . ea e . a e ca -
 . i ca i . e . . . c . . a ed i h ch . i g he . a . i ca . , i g Nai e Ba e
 . ea e . . . he . he i . ea . ed acc . ac i highe . The da a . ea . e i g af -
 f . d . . he a b i . . d h i a . e ca . . . i . . e f . a ce . . . ee i . a ce
 i . he a e a ha he . ea e a a i . i e f . ed **(htnba)** .

The e h d . . . b e f . i g a Nai e Ba e . edic i . . e . a i g i -
 . a ce , c . . a i g i . . edic i . . i h he . a . i ca . . C . . a e . . ed .
 . ea . e h . . a . i e he Nai e Ba e . edic i . ge . he . e ca . c . ec
 a c . . a ed . he . a . i ca . . Whe e f . i g a . edic i . . a e i -
 . a ce , he eaf i . . . e . . a Nai e Ba e . edic i . if i ha bee . . e
 acc . a e . e a ha he . a . i ca . . . he i e i . e . . . a . a . i ca .
 . edic i . .

T . c . . e e he e . e i e a i . , e added . i i g a d . . de . e a a i .
htnba , he e a e . e fe . ed . a **htnbp** a d **htnbaps** e ec i e .

Tab e 1 . . . a i e he c . . a . cia ed i h he ca dida e , be . d ha
 . eed ed f . **htnb** . The c . . a . cia ed i h he ada i e ch ice a e . i . .
 a fe . e . a c . . a d a i ge c . . a i . . e . . edic i . . The Nai e Ba e
 . edic i . . e . a i g i . a ce i a c . . ha ca be ha ed i h he e a a i .
 . echa i . . The c . . a . cia ed i h . a i a i g a . e a a e . edic i . . de
 a e he ge a e . e ec i e d b i g he . . age a d da e i e e . eaf . A
 . . i g i a . ch e . fe . e . . e a i . . ha a h i ge e e , highe . . i i g
 c . . . a d . . ha e . ch i ac . . he . e a . . a c . .

4 Results and Discussion

Figure 5-8 show the effect of the a1 . . . edic i . . . a egie . . . he i . . e
 a d c . . e . . ee da a , i h a d i h . . . i e .

Figure 5 show the effect of the h d f . i i g i h . . e a a e . . de .
(htnbp a d **htnbap)** d . . . e ha **ht** . A . . he Nai e Ba e . e h d . . -
 . e f . **ht** b . . gh e a a I . . d c i g . . i e i Figure 6 see **htnb**
 d i g e a , e e . . . e ha **htnbp** a d **htnbap** . The . he . e h d
 (be ide h e . i g . h . . e . . e . .) . cce f . . . e c . e he . be . i h
 i . e be . ee . he .

Table 1. Additional space/time costs beyond **htnb** requirements

	space per tree	space per leaf	time per training instance	time per test instance	time per split
htnb-stmx	cache of x instances		cache update		pass instances to leaves + NB updates
htnbp					distribution estimation
htnbps		NB model	NB update		distribution estimation
htnba		error count	NB prediction count update	decide MC or NB	
htnbap		error count	NB prediction count update	decide MC or NB	distribution estimation
htnbaps		error count NB model	NB prediction count update NB update	decide MC or NB	distribution estimation

Figure 7 illustrates the cache failure rate of the built-in leaf cache. Results are similar to the leaf cache (Figure 5). The error rate of the cache is similar to the generalised htd hash of **ht**. Once again both **htnbp** and **htnbap** perform as well as **ht**.

Adding the leaf cache to the cache of the leaf node is a good idea, but it will be seen in the next section (Figure 8) that the hash of the leaf node is a significant factor. A hash of the leaf node of 1000 instances had to be used for **htnb**, which is the same as for **ht**. Increasing the cache to 10000 instances did not improve the results.

Figure 8 demonstrates the effect of the additional **htnbp** and **htnbps** factors on the hash of the additional leaf node. The best performance is given by the leaf cache of **(htnbaps)**, which is the same as for **htnba** in fact.

The time taken to hash the additional leaf node is a significant factor in the additional UCI datasets [1] considered. The LED dataset is the most expensive, and it is a good idea to use the LED generator and test 24 bit accuracy, 10 cases, and 10% split.

The results in Figure 9 are highly interesting. The accuracy of the hash of the leaf node is 26% error, which is a significant improvement on the accuracy of **ht** (which has a accuracy of 10% error), and **htnb**. The fact that **htnb** has a higher accuracy than **ht** is a good sign. The accuracy of the hash of the leaf node is a good sign.

The results of the additional leaf node are a good sign. The accuracy of the hash of the leaf node is a good sign.

Table 2. Final accuracies achieved on tree generators

	simple tree no noise	simple tree 10% noise	complex tree no noise	complex tree 10% noise
ht	99.056 ± 0.033	82.167 ± 0.031	89.793 ± 0.168	73.644 ± 0.151
htnb	99.411 ± 0.026	81.533 ± 0.021	90.723 ± 0.153	71.425 ± 0.118
htnb-stm1k	99.407 ± 0.027	81.544 ± 0.019	90.768 ± 0.150	71.527 ± 0.108
htnb-stm10k	99.409 ± 0.025	81.593 ± 0.018	91.008 ± 0.153	71.658 ± 0.085
htnbp	97.989 ± 0.058	81.853 ± 0.042	88.326 ± 0.209	73.029 ± 0.121
htnbps	99.376 ± 0.028	82.456 ± 0.023	90.598 ± 0.153	73.063 ± 0.124
htnba	99.408 ± 0.027	82.510 ± 0.024	90.874 ± 0.153	74.089 ± 0.141
htnbap	98.033 ± 0.057	81.938 ± 0.040	88.609 ± 0.211	73.675 ± 0.127
htnbaps	99.375 ± 0.028	82.545 ± 0.024	90.935 ± 0.148	74.249 ± 0.134

Table 3. Final accuracies achieved on other datasets

	LED	Covertime
ht	72.851 ± 0.031	66.832 ± 0.163
htnb	71.645 ± 0.013	69.064 ± 0.135
htnbp	73.928 ± 0.005	68.476 ± 0.040
htnbps	73.799 ± 0.041	69.049 ± 0.145
htnba	73.935 ± 0.005	70.998 ± 0.087
htnbap	73.961 ± 0.004	71.388 ± 0.037
htnbaps	73.996 ± 0.005	71.054 ± 0.095

held the lead and held her accuracy ahead. The agent, aiming to be the best at her own task, has a 1.5% better accuracy. In the second experiment, the best performance is achieved by **htnb** with a 4.1% better accuracy.

Finally, the agent has been evaluated on the dataset generated by the File Collector domain. This collection of 581,012 instances, 10% of which are labeled, 44 binary and 7 categorical. The 10% of labels are hidden and the accuracy is measured 10 times on average. In Figure 10 we see the performance of the different methods. The best performance is achieved by **ht**. The second best is achieved by **htnb**, **htnbp** and **htnbps**. The gap between the best and the second best is very small. This is due to the fact that the best performance is achieved by **htnb** in the binary and categorical, adding the additional categorical case has a very small effect.

Table 2 and 3 show the accuracy achieved using the agent on the datasets used for the performance evaluation. Figure 5 and 10.

Overall, the results show that the best performance is achieved by the agent. In the first experiment, the agent has a 1.5% better accuracy. In the second experiment, the agent has a 4.1% better accuracy. In the third experiment, the agent has a 1.5% better accuracy. In the fourth experiment, the agent has a 1.5% better accuracy.

Overall, the results show that the best performance is achieved by the agent. In the first experiment, the agent has a 1.5% better accuracy. In the second experiment, the agent has a 4.1% better accuracy. In the third experiment, the agent has a 1.5% better accuracy. In the fourth experiment, the agent has a 1.5% better accuracy.

he di e e ce de e . . . if he e a c . . . The ada 1 e a . . . ach f htnba
 ha a e a i e . . . e head, e a i g i ca be . . . i ed . . . e ht, a de e cia
 . . . e htnb, 1 a b he . . . e e e e . . . ce-b . . . ded 1 a i . . .

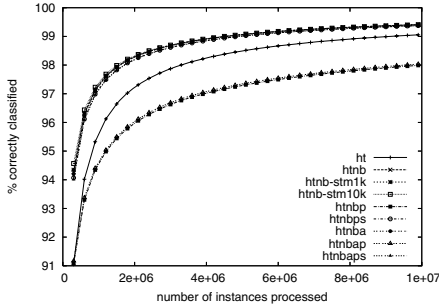


Fig. 5. Simple random tree generator with no noise

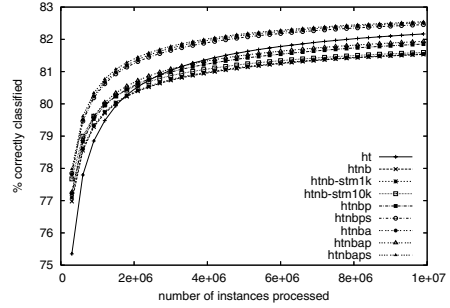


Fig. 6. Simple random tree generator with 10% noise

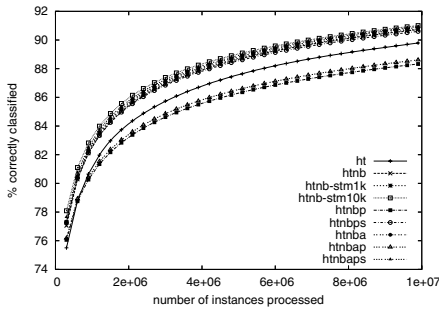


Fig. 7. Complex random tree generator with no noise

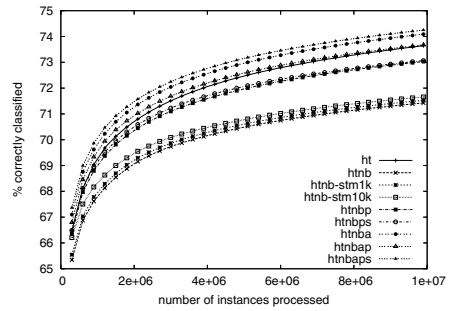


Fig. 8. Complex random tree generator with 10% noise

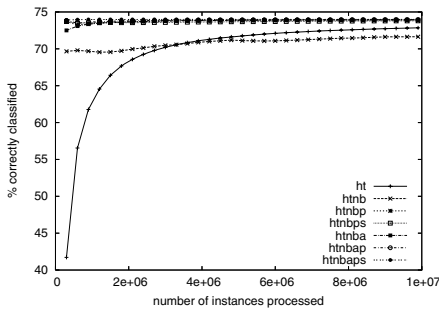


Fig. 9. LED generator

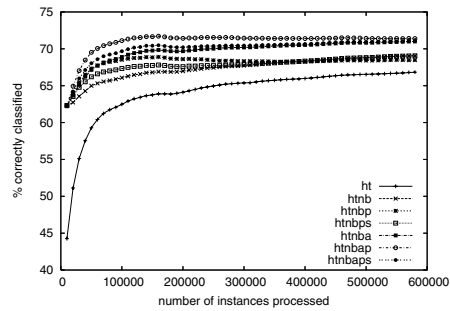


Fig. 10. Covertyp generator

5 Conclusions

Building a decision tree classifier for a large dataset is a challenging task. The performance of the classifier is highly dependent on the quality of the data and the choice of the algorithm. In this paper, we have presented a new algorithm for building decision trees. The algorithm is based on the idea of using a greedy search to find the best split at each node. The algorithm is simple and efficient, and it has been shown to perform well on a variety of datasets. The algorithm is a good choice for building decision trees in a large dataset.

References

1. C.L. Blake and C.J. Merz. UCI repository of machine learning databases, 1998.
2. Pedro Domingos and Geoff Hulten. Mining high-speed data streams. In *Knowledge Discovery and Data Mining*, pages 71–80, 2000.
3. Joao Gama, Pedro Medas, and Ricardo Rocha. Forest trees for on-line data. In *SAC '04: Proceedings of the 2004 ACM symposium on Applied computing*, pages 632–636, New York, NY, USA, 2004. ACM Press.
4. Joao Gama, Ricardo Rocha, and Pedro Medas. Accurate decision trees for mining high-speed data streams. In *KDD '03: Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 523–528, New York, NY, USA, 2003. ACM Press.

Rank Measures for Ordering

Jin Huang and Charles X. Ling

Department of Computer Science,
The University of Western Ontario,
London, Ontario, Canada N6A 5B7
{jhuang33, cling}@csd.uwo.ca

Abstract. Many data mining applications require a ranking, rather than a mere classification, of cases. Examples of these applications are widespread, including Internet search engines (ranking of pages returned) and customer relationship management (ranking of profitable customers). However, little theoretical foundation and practical guideline have been established to assess the merits of different rank measures for ordering. In this paper, we first review several general criteria to judge the merits of different single-number measures. Then we propose a novel rank measure, and compare the commonly used rank measures and our new one according to the criteria. This leads to a preference order for these rank measures. We conduct experiments on real-world datasets to confirm the preference order. The results of the paper will be very useful in evaluating and comparing rank algorithms.

1 Introduction

Ranking of cases is an increasingly important way to describe the result of many data mining and other science and engineering applications. For example, the result of document search in information retrieval and Internet search is typically a ranking of the results in the order of match. This leaves two issues to be addressed. First, given two orders of cases, how do we design or choose a measure to determine which order is better? Second, given two different rank measures of ordering, how do we tell which rank measure is more desirable?

In previous research, the issue of determining which order is better is usually addressed using accuracy and its variants, such as recall and F-measures, which are typically used in information retrieval. More recently, AUC (Area Under Curve) of the ROC (Receiver Operating Characteristics) has gained an increasing acceptance in comparing learning algorithms [1] and constructing learning models [2,3]. Bradley [4] experimentally compared popular machine learning algorithms using both accuracy and AUC, and found that AUC exhibits several desirable properties when compared to the accuracy.

However, accuracy is traditionally designed to judge the merits of classification results, and AUC is simply used as a replacement of accuracy without much reasoning for why it is a better measure, especially for the case of ordering. The main reason for this lack of understanding is that up to now, there has been no theoretical study on whether any of these measures work better than others, or whether there are even better measures in existence.

In this paper, we first review our previous work [5] on general criteria to compare two arbitrary single-number measures (see Section 2.1). Then we compare six rank measures for ordering using our general criteria. Our contributions in this part consist of a novel measure for the performance of ordering (Section 2.4), and a preference order discovered for these measures (Section 3.1). The experiments on real-world datasets confirm our analysis, which show that better rank measures are more sensitive in comparing rank algorithms (see Section 3.2).

2 Rank Measures for Ordering

In this section, we first review the criteria proposed in our previous work to compare two arbitrary measures. We then review five commonly used rank measures, and propose one new rank measure, OAUC. Then based on the comparison criteria, we will make a detailed comparison among these measures, which leads to a preference order of the six rank measures. Finally, we perform experiments with real-world data to confirm our conclusions on the preference order. The conclusions of the paper are significant for future machine learning and data mining applications involving ranking and ordering.

2.1 Review of Formal Criteria for Comparing Measures

In [5] the *degree of consistency* and *degree of discriminancy* of two measures are proposed and defined. The degree of consistency between two measures f and g , denoted as $\mathbf{C}_{f,g}$, is simply the fraction (probability) that two measures are consistent over some distribution of the instance space. Two measures are consistent when comparing two objects a and b , if f stipulates that a is better than b , g also stipulates that a is better than b . [5] define that two measures f and g are *consistent* iff the degree of consistency $\mathbf{C}_{f,g} > 0.5$. That is, f and g are consistent if they agree with each other on over half of the cases.

The *degree of discriminancy* of f over g , denoted as $\mathbf{D}_{f/g}$, is defined as the ratio of cases where f can tell the difference but g cannot, over the cases where g can tell the difference but f cannot. [5] define that a measure f is *more discriminant* (or *finer*) than g iff $D_{f/g} > 1$. That is, f is finer than g if there are more cases where f can tell the difference but g cannot, than g can tell the difference but f cannot.

2.2 Notation of Ordering

We will use some simple notations to represent ordering throughout this paper. Without loss of generality, for n examples to be ordered, we use the actual ordering position of each example as the label to represent this example in the ordered list. For example, suppose that the label of the actual highest ranked example is n , the label of the actual second highest ranked example is $n - 1$, etc. We assume the examples are ordered incrementally from left to right. Then the *true-order list* is $l = 1, 2, \dots, n$. For any ordered list generated by an ordering algorithm, it is a permutation of l . We use $\pi(l)$ to denote the ordered list generated by ordering algorithm π . $\pi(l)$ can be written as a_1, a_2, \dots, a_n , where a_i is the actual ordering position of the example that is ranked i th in $\pi(l)$.

Table 1. An example of ordered lists

l	1	2	3	4	5	6	7	8
	a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8
$\pi(l)$	3	6	8	1	4	2	5	7

Table 1 gives an instance of ordered lists with eight examples. In this table, l is the true-order list and $\pi(l)$ is the ordered list generated by an ordering algorithm π . In $\pi(l)$ from left to right are the values of a_i . We can find that $a_1 = 3, a_2 = 6, \dots, a_8 = 7$.

2.3 Previous Rank Measures for Ordering

We first review five most commonly used rank measures. Later we will invent a new rank measure which we will evaluate among the rest.

We call some of the rank measures “true-order” rank measures, because to obtain the evaluation values, we must know the true order of the original lists. Some other rank measures, however, are not true-order rank measures. They do not need the true order to obtain evaluation values; instead, only a “rough” ordering is sufficient. For example, accuracy and AUC are not true-order rank measures. As long as we know the true classification, we can calculate their values. In a sense, positive examples can be regarded as “the upper half”, and negative examples are the “lower half” in an ordering, and such a rough ordering is sufficient to obtain AUC and accuracy.

1. Euclidean Distance (ED)

If we consider the ordered list and the true order as a point (a_1, a_2, \dots, a_n) and a point $(1, 2, \dots, n)$ in an n -dimensional Euclidean space, then ED is the Euclidean Distance between these two points, which is $\sqrt{\sum_{i=1}^n (a_i - i)^2}$. For simplicity we use the squared value of Euclidean distance as the measure. Then $ED = \sum_{i=1}^n (a_i - i)^2$. Clearly, ED is a true-order rank measure.

For the example in Table 1, It is easy to obtain that $ED = (3 - 1)^2 + (6 - 2)^2 + (8 - 3)^2 + (1 - 4)^2 + (4 - 5)^2 + (2 - 6)^2 + (5 - 7)^2 + (7 - 8)^2 = 76$.

2. Manhattan Distance (MD)

This measure MD is similar to ED except that here we sum the absolute values instead of sum squared values. It is also a true-order rank measure. For our order problem $MD = \sum_{i=1}^n |a_i - i|$. For the example in Table 1, it is easy to obtain that $MD = |3 - 1| + |6 - 2| + |8 - 3| + |1 - 4| + |4 - 5| + |2 - 6| + |5 - 7| + |7 - 8| = 22$.

3. Sum of Reversed Number (SRN)

This is roughly the sum of the reversed pairs in the list. That is, $SRN = \sum_{i=1}^n s(i)$. It is clearly a true-order measure.

For the i th example, its reversed number $s(i)$ is defined as the number of examples whose positions in $\pi(l)$ are greater than i but the actual ranked positions are less than i . For the example in Table 1, we can find that the examples of 1 and 2 are both ranked higher than the first example 3 in $\pi(l)$. Thus $s(1) = 1 + 1 = 2$. Similarly we have $s(2) = 4, s(3) = 5$, etc. Therefore the SRN for the ordered list $\pi(l)$ is $SRN = 2 + 4 + 5 + 0 + 1 + 0 + 0 + 0 = 12$.

4. **Area Under Curve (AUC)**

The Area Under the ROC Curve, or simply AUC, is a single-number measure widely used in evaluating classification algorithms, and it is not a true-order measure for ranking. To calculate AUC for an ordered list, we only need the true classification (positive or negative examples). For a balanced ordered ranked list with n examples (half positive and half negative), we treat any example whose actual ranked position is greater than $\frac{n}{2}$ as a positive example; and the rest as negative. From left to right we assume the ranking positions of positive examples are $r_1, r_2, \dots, r_{\lfloor \frac{n}{2} \rfloor}$.

Then $AUC = \frac{\sum_{a_{r_i} > n/2} (r_i - i)}{n^2}$ [6].

In Table 1, 5, 6, 7, and 8 are positive examples positioned at 2, 3, 7, and 8 respectively. Thus, $AUC = \frac{(2-1)+(3-2)+(7-3)+(8-4)}{4 \times 4} = \frac{5}{8}$.

5. **Accuracy (acc)**

Like AUC, accuracy is also not a true-order rank measure. Similar to AUC, if we classify examples whose rank position above half of the examples as positive, and the rest as negative, we can calculate accuracy easily as $acc = \frac{tp+tn}{n}$, where tp and tn are the number of correctly classified positive and negative examples respectively. In the ordered list $\pi(l)$ in Table 1, 5, 6, 7, and 8 are positive examples, others are negative examples. Thus $tp = 2, tn = 2$. $acc = \frac{2+2}{8} = \frac{1}{2}$.

2.4 **New Rank Measure for Ordering**

We propose a new measure called Ordered Area Under Curve (OAUC), as it is similar to AUC both in meaning and calculation. The only difference is that each term in the formula is weighted by its true order, and the sum is then normalized. Thus, OAUC is a true-order measure. This measure is expected to be better than AUC since it ‘‘spreads’’ its values more widely compared to AUC.

OAUC is defined as follows:

$$OAUC = \frac{\sum a_{r_i} (r_i - i)}{\lfloor \frac{n}{2} \rfloor \sum_{i=1}^{\lfloor \frac{n}{2} \rfloor} (\lfloor \frac{n}{2} \rfloor + i)}$$

In the ordered list in Table 1, the positive examples are 5, 6, 7, 8 which are positioned at 7, 2, 8 and 3 respectively. Thus $r_1 = 2, r_2 = 3, r_3 = 7, r_4 = 8$, and $a_{r_1} = 6, a_{r_2} = 8, a_{r_3} = 5, a_{r_4} = 7$. $OAUC = \frac{6(2-1)+8(3-2)+5(7-3)+7(8-4)}{4((4+1)+(4+2)+(4+3)+(4+4))} = \frac{31}{52}$.

3 **Comparing Rank Measures for Ordering**

We first intuitively compare some pairs of measures and analyze whether any two measures satisfy the criteria of consistency and discriminancy. To begin with, we consider ED and MD because these two measures are quite similar in their definitions except that ED sums the squared distance while MD sums the absolute value. We expect that these two measures are consistent in most cases. On the other hand, given a dataset with n examples there are a total of $O(n^3)$ different ED values and $O(n^2)$ different MD values. Thus ED is expected to be more discriminant than MD. Therefore we expect that ED is consistent with and more discriminant than MD.

For AUC and OAUC, since OAUC is an extension of AUC, intuitively we expect that they are consistent. Assuming there are n_1 negative examples and n_0 positive examples, the different values for OAUC is $n_1 \sum_{i=1}^{n_0} (n_1 + i)$, which is greater than the different values of AUC ($n_0 n_1$). We can also expect that OAUC is more discriminant and therefore better than AUC.

However for the rest of the ordering measures we cannot make these intuitive claims because they have totally different definitions or computational methods. Therefore, in order to perform an accurate and detailed comparison and to verify or overturn our intuitions, we will conduct experiments to compare all measures.

3.1 Comparing Rank Measures on Artificial Datasets

To obtain the average degrees of consistency and discriminancy for all possible ranked lists, we use artificial datasets which consist of all possible ordered list of length 8.¹ We assume that the ordered lists are uniformly distributed. We exhaustively compare all pairs of ordered lists and calculate the degree of consistency and degree of discriminancy between two rank measures for ordering.

Table 2 lists the degree of consistency between every pair of six rank measures for ordering. The number in each cell represents the degree of consistency between the measures in the same row and column of the cell. We can find that the degree of consistency between any two measures are greater than 0.5, which indicates that these measures are “similar” in the sense that they are more likely to be consistent than inconsistent.

Table 3 shows the degree of discriminancy among all 6 rank measures. The number in the cell of the i th row and the j th column is the degree of discriminancy for the measure in i th row over the one in j th column.

From these two tables we can draw the following conclusions. First, these results verified our previous intuitive conclusions about the relations between ED and MD, and between AUC and OAUC. The degree of consistency between ED and MD is 0.95, and between AUC and OAUC 0.99, which means that ED and MD, and AUC and OAUC are highly consistent. The degree of discriminancy for ED over MD, and for OAUC over AUC are greater than 1, which means that ED is better than MD, and OAUC is better than AUC.

Table 2. Degree of consistency between pairs of rank measures for ordering

	AUC	SRN	MD	ED	OAUC	acc
AUC	1	0.88	0.89	0.87	0.99	0.98
SRN	0.88	1	0.95	0.98	0.89	0.91
MD	0.89	0.95	1	0.95	0.90	0.95
ED	0.87	0.98	0.95	1	0.88	0.90
OAUC	0.99	0.89	0.90	0.88	1	0.97
acc	0.98	0.91	0.95	0.90	0.97	1

¹ There are $n!$ different ordered lists for length n , so it is infeasible to enumerate longer lists.

Table 3. Degree of discriminancy between pairs of rank measures for ordering

	AUC	SRN	MD	ED	OAUC	acc
AUC	1	0.88	1.42	0.21	0.0732	14.0
SRN	1.14	1	1.84	0.242	0.215	9.94
MD	0.704	0.54	1	0.117	0.116	6.8
ED	4.76	4.13	8.55	1	0.87	38.2
OAUC	13.67	4.65	8.64	1.15	1	94.75
acc	0.071	0.10	0.147	0.026	0.011	1

Second, since all values of the degree of consistency among all measures are greater than 0.5, we can decide which measure is better than another only based on the value of degree of discriminancy. Recall (Section 2.1) that a measure f is better than another measure g iff $C_{f,g} > 0.5$ and $D_{f/g} > 1$. The best measure should be the one whose degrees of discriminancy over all other measures are greater than 1. From Table 3 we can find that all the numbers in the OAUC row are greater than 1, which means that the measure OAUC's degrees of discriminancy over all other measures are greater than 1. Therefore OAUC is the best measure. In the same way we can find that ED is the second best measure, and SRN is the third best. The next are AUC, MD, and acc is the worst.

Finally we can obtain the following preference order of for all six rank measures for ordering:

$$OAUC \succ ED \succ SRN \succ AUC \succ MD \succ acc$$

From the preference order we can conclude that OAUC, a new measure we design based on AUC, is the best measure. ED is the close, second best. The difference for these two measures are not very large (the degree of discriminancy for OAUC over ED is only 1.15). Therefore we should use OAUC and ED instead of others to evaluate ordering algorithms in most cases. Further, the two non-true-order classification measures AUC and accuracy do not perform well as compared with the true-order measures ED and SRN. This suggests that generally we should avoid using classification measures such as AUC and accuracy to evaluate ordering. Finally, MD is the worst true-order measure, and it is even worse than AUC. It should be avoided.

3.2 Comparing Rank Measures with Ranking Algorithms

In this section, we perform experiments to compare two classification algorithms in terms of the six rank measures. What we hope to conclude is that the better rank measures (such as OAUC and ED) would be more sensitive to the significance test (such as the t-test) than other less discriminant measures (such as MD and accuracy). That is, OAUC and ED are more likely to tell the difference between two algorithms than MD and accuracy can. Note that here we do not care about which rank algorithm predicts better; we only care about the sensitivity of the rank measures that are used to compare the rank algorithms. The better the rank measure (according to our criteria), the more sensitive it would be in the comparison, and the more meaningful the conclusion would be for the comparison.

We choose Artificial Neural Networks (ANN) and Instance-Based Learning algorithm (IBL) as our algorithms as they can both accept and produce continuous target. The ANN that we use has one hidden layer; the number of nodes in the hidden layer is half of the input layer (the number of attributes). We use real-world datasets to evaluate and compare ANN and IBL with the six rank measures. We select three real-world datasets *Wine*, *Auto-Mpg* and *CPU-Performance* from the UCI Machine Learning Repository [7].

In our experiments, we run ANN and IBL with the 10-fold cross validation on the training datasets. For each round of the 10-fold cross validation we train the two algorithms on the same training data and test them on the same testing data. We measure the testing data with six different rank measures (OAUC, ED, SRN, AUC, MD and acc) discussed earlier in the paper. We then perform paired, two-tailed t-tests on the 10 testing datasets for each measure to compare these two algorithms.

Table 4 shows the significance level in the t-test.² The smaller the values in the table, the more likely that the two algorithms (ANN and IBL) are significantly different, and the more sensitive the measure is when it is used to compare the two algorithms. Normally a threshold is set up and a binary conclusion (significantly different or not) is obtained. For example, if we set the threshold to be 0.95, then for the artificial dataset, we would conclude that ANN and IBL are statistically significantly different in terms of ED, OAUC and SRN, but not in terms of AUC, MD and acc. However, the actual significance level in Table 4 is more discriminant for the comparison. That is, it is “a better measure” than the simple binary classification of being significantly different or not.

Table 4. The significance level in the paired t-test when comparing ANN and IBL using different rank measures

Measures	Wine	Auto-mpg	CPU
OAUC	0.031	8.64×10^{-4}	1.48×10^{-3}
ED	0.024	1.55×10^{-3}	4.01×10^{-3}
SRN	0.053	8.89×10^{-3}	5.91×10^{-3}
AUC	0.062	5.77×10^{-3}	8.05×10^{-3}
MD	0.053	0.0167	5.97×10^{-3}
acc	0.126	0.0399	0.0269

From Table 4 we can obtain the preference order from the most sensitive measure (the smallest significance level) to the least sensitive measure (the largest significance level) for each dataset is:

- Wine: ED, OAUC, SRN = MD, AUC, acc.
- Auto-mpg: OAUC, ED, AUC, SRN, MD, acc.
- CPU-Performance: OAUC, ED, SRN, MD, AUC, acc.

These preference orders are roughly the same as the preference order of these measures discovered in the last section:

² The confidence level for the two arrays of data to be statistically different is one minus the values in the table.

$$OAUC \succ ED \succ SRN \succ AUC \succ MD \succ acc$$

The experimental results confirm our analysis in the last section. That is, OAUC and ED are the best rank measures for evaluating orders. In addition, MD and accuracy should be avoided as rank measures. These conclusions will be very useful for comparing and constructing machine learning algorithms for ranking, and for applications such as Internet search engines and data mining for CRM (Customer Relationship Management).

4 Conclusions

In this paper we use the criteria proposed in our previous work to compare five commonly used rank measures for ordering and a new proposed rank measure (OAUC). We conclude that OAUC is actually the best rank measure for ordering, and it is closely followed by the Euclidian distance (ED). Our results indicate that in comparing different algorithms for the order performance, we should use OAUC or ED, and avoid the least sensitive measures such as Manhattan distance (MD) and accuracy.

In our further work, we plan to improve existing rank learning algorithms by optimizing the better measures, such as OAUC and ED, discovered in this paper.

References

1. Provost, F., Domingos, P.: Tree induction for probability-based ranking. *Machine Learning* **52:3** (2003) 199–215
2. Ferri, C., Flach, P.A., Hernandez-Orallo, J.: Learning decision trees using the area under the ROC curve. In: *Proceedings of the Nineteenth International Conference on Machine Learning (ICML 2002)*. (2002) 139–146
3. Ling, C.X., Zhang, H.: Toward Bayesian classifiers with accurate probabilities. In: *Proceedings of the Sixth Pacific-Asia Conference on KDD*. Springer (2002) 123–134
4. Bradley, A.P.: The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognition* **30** (1997) 1145–1159
5. Ling, C.X., Huang, J., Zhang, H.: AUC: a statistically consistent and more discriminating measure than accuracy. In: *Proceedings of 18th International Conference on Artificial Intelligence (IJCAI-2003)*. (2003) 519–526
6. Hand, D.J., Till, R.J.: A simple generalisation of the area under the ROC curve for multiple class classification problems. *Machine Learning* **45** (2001) 171–186
7. Blake, C., Merz, C.: *UCI Repository of machine learning databases* [<http://www.ics.uci.edu/~mllearn/MLRepository.html>]. Irvine, CA: University of California, Department of Information and Computer Science (1998)

Dynamic Ensemble Re-Construction for Better Ranking

Jin Huang and Charles X. Ling

Department of Computer Science,
The University of Western Ontario,
London, Ontario, Canada N6A 5B7
{jhuang33, cling}@csd.uwo.ca

Abstract. Ensemble learning has been shown to be very successful in data mining. However most work on ensemble learning concerns the task of classification. Little work has been done to construct ensembles that aim to improve ranking. In this paper, we propose an approach to re-construct new ensembles based on a given ensemble with the purpose to improve the ranking performance, which is crucial in many data mining tasks. The experiments with real-world data sets show that our new approach achieves significant improvements in ranking over the original Bagging and Adaboost ensembles.

1 Introduction

Classification is one of the fundamental tasks in knowledge discovery and data mining. The performance of a classifier is usually evaluated by predictive accuracy. However, most machine learning classifiers can also produce the probability estimation of the class prediction. Unfortunately, this probability information is ignored in the measure of accuracy.

In many real-world data mining applications, however, we often need the probability estimations or ranking. For example, in direct marketing, we often need to promote the most likely customers, or we need to deploy different promotion strategies to customers according to their likelihood of purchasing. To accomplish these tasks we need a ranking of customers according to their likelihood of purchasing. Thus ranking is often more desirable than classification in these data mining tasks.

One natural question is how to evaluate a classifier's ranking performance. In recent years, the area under the ROC (Receiver Operating Characteristics) curve, or simply AUC, is increasingly received attention in the communities of machine learning and data mining. Data mining researchers [1,2] have shown that AUC is a good summary in measuring a classifier's overall ranking performance. Hand and Till [3] present a simple approach to calculating AUC of a classifier for binary classification.

$$\hat{A} = \frac{S_0 - n_0(n_0 + 1)/2}{n_0 n_1}, \quad (1)$$

where n_0 and n_1 are the numbers of positive and negative examples respectively, and $S_0 = \sum r_i$, where r_i is the rank of the i_{th} positive example in the ranked list.

Ensemble is a general approach which trains a number of classifiers and then combines their predictions in classification. Many researches [4,5,6] have shown that the

ensemble is quite effective in improving the classification accuracy compared with a single classifier. The reason is that the prediction error of an individual classifier can be counteracted by the combination with other classifiers. Bagging [5] and Boosting [7] are two of the most popular ensemble techniques.

Most previous work of ensemble learning is focussed on classification. To our knowledge, there is little work that directly constructs ensembles to improve probability estimations or ranking. [8] compared the probability estimations (ranking) performance of different learning algorithms by using AUC as the comparison measure and demonstrated that Boosted trees and Bagged trees perform better in terms of ranking than Neural Networks and SVMs. [9] used the boosting technique on the general preference learning (ranking) problem and proposed a new ranking boosting algorithm: RankBoost.

In this paper, we propose a novel approach to improve the ranking performance over a given ensemble. The goal of this approach is to select some classifiers from the given ensemble to re-construct new ensembles. It first uses the k -Nearest Neighbor method to find training data subsets which are most similar to the test set, then it uses the measure SAUC (see Section 2.2) as heuristic to dynamically choose the diverse and well performed classifiers. This approach is called DERC (Dynamic Ensemble Re-Construction) algorithm. The new ensembles constructed by this approach are expected to have better ranking performance than the original ensemble.

The paper is organized as follows. In Section 2 we give detailed description for our new algorithm. In Section 3 we perform experiments on real world data sets to show the advantages of the new algorithm.

2 DERC (Dynamic Ensemble Re-Construction) Algorithm

In an ensemble, the combination of the predictions of several classifiers is only useful if they disagree to some degree. Each ensemble classifier may perform diversely during classification. Our DERC algorithm is motivated by this diversity property of ensemble. The diversity implies that each ensemble classifier performs best in probability estimation (ranking) only in a subset of training instances. Thus given a test (sub)set, if we use the k -Nearest Neighbor method to find some training subsets that are most similar to it, the classifiers that perform diversely and accurately on those similar training subsets are also expected to perform well on the test (sub)set. Therefore the new ensembles constructed are expected to have better ranking performance than the original ensemble.

Our DERC algorithm involves two basic steps: finding the most similar training (sub)sets, and selecting the diverse and accurate classifiers.

Now we use Figure 1 to illustrate how DERC algorithm works. Suppose that we are given an ensemble E with multiple classifiers built on a training set S , and we have an unlabeled test set T at hand. Our goal is to select some classifiers from the ensemble E to build one or more new ensembles to perform ranking on test set T .

2.1 Finding the Most Similar Training Subsets

The first step is to stratify the test set to some equal parts and find the most similar training subsets corresponding to test partitions. Since the labels of test instances are

unknown, we randomly pick a classifier from ensemble E to classify the test set T to obtain the predicted labels. Assume that we want to construct 3 new ensembles. According to the predicted class labels we stratify (partition with equal class distributions) the test set T into 3 equal sized parts: T_1 , T_2 , and T_3 . We want to select some classifiers from ensemble E to build 3 different new ensembles which are responsible for ranking T_1 , T_2 and T_3 respectively.

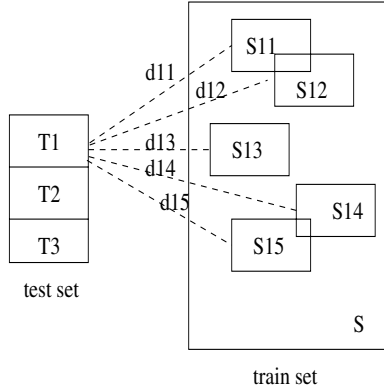


Fig. 1. An example for the similar sets

For each stratified test subset we use the k -Nearest Neighbor method to find k subsets of training set which are most similar to that test set part. For each instance of the test subset, we compute the distances from this instance to all training instances and find the nearest k instances. We use the following method to compute the distance between two instances u and v , which are from the test subset and training dataset, respectively. Suppose that an instance has k_1 nominal attributes A_i and k_2 numerical attributes B_j . We use the simplified VDM measure proposed in [10] to compute the distance of all nominal attributes.

$$VDM(u, v) = \sum_C \sum_{i=1}^{k_1} \left(\frac{N_{A_i=a_u, C=c}}{N_{A_i=a_u}} - \frac{N_{A_i=a_v, C=c}}{N_{A_i=a_v}} \right)^2$$

where $N_{A_i=a_u}$ is the number of instances in test subset holding value a_u on attribute A_i , $N_{A_i=a_u, C=c}$ is the number of instances in test subset which are predicted belonging to class c and hold value a_u on attribute A_i . Here note that since test set is unlabeled, we use the class labels predicted in the first step.

We simply use the Euclidean distance to compute the difference of numerical attributes. $ED(u, v) = \sum_{i=1}^{k_2} (b_{u_i} - b_{v_i})^2$, where b_{u_i} is instance u ' value on numerical attribute B_i .

The distance of u and v is

$$d(u, v) = VDM(u, v) + ED(u, v)$$

After the distances are computed, we randomly pick one from the k nearest instances of each test instance and use them to form a training subset. This subset is most similar to the test subset. We can use this method to find a desired number of most similar training subsets. The distance between two similar data sets is simply the average distances of each test subset instance with its corresponding nearest training instance. As shown in Figure 1, assume that $S_{11}, S_{12}, S_{13}, S_{14}$ and S_{15} are T_1 's 5 most similar training subsets. Their distances to T_1 are computed as $d_{11}, d_{12}, d_{13}, d_{14}$ and d_{15} , respectively.

2.2 Selecting Diverse and Accurate Classifiers

After the most similar training subsets are found, we use the following strategy to select diverse and accurate classifiers from original ensemble. Instead of directly using AUC as the criterion to choose classifiers, we propose a new measure SAUC (Softened Area Under the ROC Curve) as the heuristic.

For a binary classification task, SAUC is defined as

$$SAUC(\gamma) = \frac{\sum_{i=1}^m \sum_{j=1}^n U(p_i^+ - p_j^-)^\gamma}{mn} \tag{2}$$

$$U(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

where $\gamma \geq 0$, p_i^+ , p_j^- represent the predicted probabilities of being positive for the i th positive and the j th negative examples in all m positive examples and n negative examples, respectively.

We choose a series of measures $SAUC(\gamma_1), SAUC(\gamma_2), \dots, SAUC(\gamma_n)$ as heuristics. We use SAUCs as heuristics for two reasons. First, SAUC with different powers γ may have different sensitivities and robustness to instance ranking variations. Thus using SAUCs with varied power γ as heuristics can more reliably select diverse classifiers in terms of ranking. Second, SAUC is a softened version of AUC and thus it is basically consistent with AUC. From Equation 2 we can see that $SAUC(0) = AUC$. Thus using SAUCs as criteria can select the classifiers with accurate ranking performance.

As shown in Figure 1 we use each classifier C_t of ensemble E to classify $S_{11}, S_{12}, S_{13}, S_{14}$ and S_{15} to obtain the respective $SAUC(\gamma_1)$ as $SA_{11}, SA_{12}, SA_{13}, SA_{14}, SA_{15}$. We then compute a score for C_t , which is the weighted average of the $SAUC(\gamma_1)$ values obtained above. It is $S_t = \sum_{i=1}^5 \frac{SA_{1i}}{d_{1i}}$. We choose the classifier with the highest score. We repeat the above step n times by using a different $SAUC(\gamma_i)$ each time to select a new classifier.

Finally we use all the classifiers selected to construct a new ensemble. This ensemble is responsible for ranking T_1 . The new ensemble combination method is weighted averaging, in which a classifier's weight is its score computed above. Using the same method we can construct two other ensembles which are responsible for ranking T_2 and T_3 , respectively. We give the pseudo-code of this algorithm in Table 1.

One natural question about the DERC algorithm is that how many new ensembles should be constructed to give the best ranking performance. Since the number of test set partitions equals to the number of new ensembles, this question is equivalent to

Table 1. The pseudo code for DERC algorithm

```

DERC( $E, S, T, n$ )
  Input:
   $E$  : An ensemble with classifiers  $C_1, \dots, C_N$ 
   $S$  : Training data set
   $T$  : Test data set
   $n$  : The number of test set partitions

  choose a classifier from  $E$  to classify  $T$ 
  stratify  $T$  into  $T_1, T_2, \dots, T_n$ 
  for each partition  $T_i$  do
     $E_i^* \leftarrow \phi$ 
    find the most similar training subsets  $S_{i1}, S_{i2}, \dots, S_{ik}$ 
    compute the distances  $d_{i1}, d_{i2}, \dots, d_{ik}$  from  $T_i$  to  $S_{i1}, \dots, S_{ik}$  respectively
    for each measure  $SAUC(\gamma_u)$  do
      for each classifier  $C_t$  do
        run  $C_t$  on  $S_{i1}, S_{i2}, \dots, S_{ik}$ 
        obtain the  $SAUC(\gamma_u)$  of  $C_t$  as  $SA_{t_{i1}}, \dots, SA_{t_{ik}}$ 
        compute the ranking score for classifier  $C_t$ 
         $r_t \leftarrow \sum_{j=1}^k \frac{SA_{t_{ij}}}{d_{ij}}$ 
      endfor
      choose the classifier  $CC$  with highest score  $r_t$ 
       $E_i^* \leftarrow E_i^* \cup CC$ 
    endfor
  endfor
  return all  $E_i^*$ 

```

how to choose an optimal number of test set partitions. Clearly, a small number of partitions generally means large partitioned test subsets, which corresponds to large similar training subsets. Thus the corresponding new ensemble may not specialize on all instances of the similar training subsets. Therefore our algorithm may not perform best on a small number of partitions. On the contrary for very large number of partitions, the size of similar training subsets will be very small. In this case there is a danger of overfitting. Therefore we can claim that generally too small or too large number of partitions should be avoided. We will perform experiments in the next section to confirm this claim.

3 Experimental Evaluation

To evaluate the performance of our algorithm, we extract 16 representative binary data sets from UCI [11].

We use Bagging and Adaboost as the ensembling methods and Naive Bayes as the base learner. We choose WEKA [12] as the implementations. In order to increase the ensemble diversity, we randomly select half of the training data for each bootstrap in our

Bagging process. This can guarantee that the bagging classifiers are diverse to some degree. We compare the performance of DERC with Bagging and Adaboost respectively.

In our DERC algorithm we use $SAUC(\gamma_i)$ as criteria to select classifiers. We have to determine the suitable number and scores of the powers γ_i by taking into account the tradeoff between the quality of results and computational costs. We test the SAUC with a wide ranges of powers γ by using all the 16 datasets in the our experiments. The analysis of these measures' performance shows that the power range of [0,3] is a good choice for SAUC. We choose 9 different SAUC with the powers of 0, 0.1, 0.4, 0.8, 1.0, 1.5, 2, 2.5, 3 in our experiment.

We follow the procedure below to perform our experiment:

1. We discretize the continuous attributes in all data sets using the entropy-based method described in [13].
2. We perform 5-fold cross validation on each data set. In each fold we train an ensemble with 15 classifiers using Bagging and Adaboost methods, respectively. We then run our DERC algorithm on the ensemble trained. By varying the number of test set partitions, we have a number of different DERC algorithm models.
3. We run the second step 20 times and we compute the average AUC for all the predictions.

We use a common statistic to compare the learning algorithms across all data sets. We performed two tailed paired t-test with 95% confidence level to count in how many datasets one algorithm performs significantly better, same, and worse than another algorithm respectively. We use win/draw/loss to represent this.

The experimental results are listed in Table 2 and Table 3.

Table 2. Comparing the predictive AUC of DERC algorithms with Bagging

Dataset	Bagging	DERC(1)	DERC(2)	DERC(3)	DERC(4)	DERC(6)
breast	98.84 ± 0.56	98.84 ± 0.53	98.83 ± 0.50	98.85 ± 0.59	98.86 ± 0.55	98.81 ± 0.59
cars	93.56 ± 3.0	94.77 ± 2.2	94.9 ± 2.7	94.83 ± 2.7	94.87 ± 2.9	95.02 ± 2.1
credit	92.89 ± 1.2	93.43 ± 1.1	93.36 ± 1.2	93.32 ± 1.2	93.3 ± 1.4	93.3 ± 1.1
echocardio	72.34 ± 8.4	72.34 ± 8.4	74.21 ± 8.3	74.11 ± 8.4	74.11 ± 8.4	73.09 ± 8.4
eco	99.28 ± 0.84	99.34 ± 1.1	99.34 ± 1.0	99.32 ± 0.84	99.3 ± 1.0	99.33 ± 0.84
heart	85.89 ± 0.45	86.01 ± 0.5	86.97 ± 0.5	86.81 ± 0.64	86.06 ± 1.7	85.92 ± 2.6
hepatitis	86.73 ± 2.6	87.06 ± 2.6	87.5 ± 2.9	89.14 ± 2.6	88.59 ± 2.4	88.2 ± 1.8
import	97.75 ± 2.6	97.75 ± 2.6	97.59 ± 2.8	97.72 ± 2.6	97.72 ± 2.6	97.74 ± 2.6
liver	61.77 ± 1.6	61.33 ± 0.45	61.64 ± 0.6	61.4 ± 0.18	61.26 ± 0.3	61.19 ± 3.7
pima	77.27 ± 8.9	79.33 ± 7.6	79.29 ± 7.7	79.26 ± 8.0	79.14 ± 8.6	79.22 ± 8.7
thyroid	95.12 ± 1.7	95.19 ± 1.6	95.10 ± 1.6	95.16 ± 1.9	95.24 ± 1.9	95.29 ± 1.5
voting	96.00 ± 0.36	96.08 ± 0.36	96.07 ± 0.36	96.27 ± 0.36	95.99 ± 0.36	96.01 ± 0.36
sick	96.84 ± 1.56	95.20 ± 2.48	●94.27 ± 2.11	●94.27 ± 3.47	●93.99 ± 2.79	●94.08 ± 3.02
ionosphere	94.59 ± 3.21	94.80 ± 3.22	95.96 ± 3.47	95.85 ± 2.63	95.84 ± 2.79	95.84 ± 3.92
german	84.26 ± 4.02	87.58 ± 4.33	87.40 ± 4.1	87.23 ± 4.21	87.44 ± 4.2	87.4 ± 4.17
mushroom	99.89 ± 0.04	99.79 ± 0.04	99.88 ± 0.04	99.90 ± 0.04	99.89 ± 0.04	99.89 ± 0.04
w/d/l		4/12/0	7/8/1	8/7/1	8/7/1	7/8/1

Table 3. Comparing the predictive AUC of DERC algorithms with Adaboost

Dataset	AdaBoost	DERC(1)	DERC(2)	DERC(3)	DERC(4)	DERC(6)
breast	98.99 ± 2.1	98.39 ± 2.4	98.41 ± 2.1	98.46 ± 2.1	98.51 ± 2.1	98.53 ± 2.1
cars	91.74 ± 5.0	93.21 ± 5.0	93.14 ± 5.0	94.72 ± 5.0	93.89 ± 5.0	93.21 ± 5.0
credit	92.06 ± 3.7	92.04 ± 3.7	92.06 ± 3.5	92.08 ± 4.8	92.10 ± 5.3	91.77 ± 4.7
echocardio	72.02 ± 4.8	73.94 ± 4.8	73.94 ± 4.8	73.94 ± 4.8	73.94 ± 4.8	73.94 ± 4.8
eco	99.30 ± 1.0	99.13 ± 1.0	99.02 ± 1.0	99.24 ± 1.0	99.27 ± 1.0	99.62 ± 1.0
heart	88.03 ± 0.28	88.51 ± 0.31	90.39 ± 0.28	89.72 ± 1.22	89.34 ± 1.6	89.58 ± 1.24
hepatitis	85.25 ± 8.6	●83.16 ± 5.8	●83.03 ± 5.6	●83.24 ± 8.6	●83.9 ± 8.8	●83.84 ± 5.4
import	98.99 ± 1.7	98.90 ± 0.0	98.98 ± 3.6	98.73 ± 0.0	98.68 ± 5.2	98.88 ± 0.0
liver	65.45 ± 6.2	66.44 ± 4.1	66.20 ± 5.1	67.08 ± 5.1	67.77 ± 5.1	66.29 ± 5.1
pima	75.99 ± 8.3	74.92 ± 8.1	74.89 ± 7.2	77.81 ± 8.3	77.99 ± 6.5	78.13 ± 8.4
thyroid	95.61 ± 0.35	95.55 ± 0.8	95.64 ± 0.27	95.65 ± 0.18	95.58 ± 0.71	95.58 ± 0.35
voting	96.37 ± 2.9	96.32 ± 2.9	96.39 ± 3.3	96.5 ± 1.4	96.37 ± 1.4	96.37 ± 2.9
sick	97.02 ± 1.56	97.08 ± 1.51	97.07 ± 1.43	97.02 ± 1.5	96.99 ± 1.24	97.08 ± 2.54
ionosphere	94.56 ± 3.21	94.80 ± 3.47	95.96 ± 4.37	95.85 ± 4.26	95.84 ± 3.97	95.84 ± 3.68
german	86.41 ± 4.02	88.24 ± 4.33	88.21 ± 4.1	88.21 ± 4.21	88.19 ± 4.2	88.19 ± 4.17
mushroom	99.92 ± 0.04	99.79 ± 0.04	99.88 ± 0.04	99.90 ± 0.04	99.89 ± 0.04	99.89 ± 0.04
w/d/l		3/12/1	4/11/1	5/10/1	5/10/1	4/11/1

Table 2 shows the AUC values for the Bagging algorithm and the DERC algorithms with different settings on various data sets. We use DERC(i) to denote the corresponding DERC algorithm which generate a number of i new ensembles. Each data cell represents the average AUC value of the 20 trials of 5-fold cross validation for the corresponding algorithm and data set. The data in bold shows the corresponding algorithm performs significantly better than Bagging on the corresponding data set. The data with a “●” means it is significantly worse than that of Bagging.

From this table, we can see that DERC outperforms the original Bagging algorithm. The w/d/l statistics shows that all DERCs with different settings have much more wins than losses compared with Bagging algorithm. If we rank them according to the w/d/l number, we can see that the DERC with 3 or 4 partitions performs best, the DERC with 2 or 6 partitions the second best, while the DERC with 1 partition the worst.

We can also see how the partition numbers influences the dynamic re-construction performance. We can observe that generally the dynamic re-constructions with the partition numbers of 3 or 4 perform best. It shows that dynamic re-construction with intermediate number of partitions outperforms that with large or small number of partitions. This result confirms our discussion in the previous section.

We also compare our DERC algorithm with Adaboost and report the results in Table 3. The similar comparisons show that DERC also significantly outperforms Adaboost in terms of AUC. DERC(3) wins in 5 datasets, ties in 10 datasets on loses only in 1 dataset.

4 Conclusions and Future Work

In this paper we propose a novel dynamic re-construction technique which aims to improve the ranking performance of any given ensemble. This is a generic technique

which can be applied on any existing ensembles. The advantage is that it is independent of the specific ensemble construction method. The empirical experiments show that this dynamic re-construction technique can achieve significant performance improvement in term of ranking over the original Bagging and Adaboost ensembles, especially with an intermediate number of partitions .

In our current study we use Naive Bayes as the base learner. For our future work, we plan to investigate how other learning algorithms perform with the DERC technique. We also plan to explore whether DERC is also effective when it is applied on other ensemble methods.

References

1. Bradley, A.P.: The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognition* **30** (1997) 1145–1159
2. Provost, F., Fawcett, T.: Analysis and visualization of classifier performance: comparison under imprecise class and cost distribution. In: *Proceedings of the Third International Conference on Knowledge Discovery and Data Mining*. AAAI Press (1997) 43–48
3. Hand, D.J., Till, R.J.: A simple generalisation of the area under the ROC curve for multiple class classification problems. *Machine Learning* **45** (2001) 171–186
4. Bauer, E., Kohavi, R.: An empirical comparison of voting classification algorithms: Bagging, boosting and variants. *Machine Learning* **36** (1999) 105–139
5. Breiman, L.: Bagging predictors. *Machine Learning* **24** (1996) 123–140
6. Quinlan, J.R.: Bagging, boosting, and C4.5. In: *Proceedings of the Thirteenth National Conference on Artificial Intelligence*. (1996) 725 – 730
7. Freund, Y., Schapire, R.E.: Experiments with a new boosting algorithm. In: *Proceedings of the Thirteenth International Conference on Machine Learning*. (1996) 148 – 156
8. Caruana, R., Niculescu-Mizil, A.: An empirical evaluation of supervised learning for ROC area. In: *The First Workshop on ROC Analysis in AI*. (2004)
9. Freund, Y., Iyer, R., Schapire, R., Singer, Y.: An efficient boosting algorithm for combining preferences. *Journal of Machine Learning Research* **4** (2003) 933–969
10. Stanfill, C., Waltz, D.: Toward memory-based reasoning. *Communications of the ACM* **29** (1986) 1213–1228
11. Blake, C., Merz, C.: UCI repository of machine learning databases. <http://www.ics.uci.edu/~mlern/MLRepository.html> (1998) University of California, Irvine, Dept. of Information and Computer Sciences.
12. Witten, I.H., Frank, E.: *Data Mining: Practical machine learning tools with Java implementations*. Morgan Kaufmann, San Francisco (2000)
13. Fayyad, U., Irani, K.: Multi-interval discretization of continuous-valued attributes for classification learning. In: *Proceedings of Thirteenth International Joint Conference on Artificial Intelligence*. Morgan Kaufmann (1993) 1022–1027

Frequency-Based Separation of Climate Signals

Alexander Ilin¹ and Harri Valpola²

¹ Helsinki University of Technology, Neural Networks Research Centre,
P.O. Box 5400, FI-02015 TKK, Espoo, Finland
Alexander.Ilin@tkk.fi

² Helsinki University of Technology, Lab. of Computational Engineering,
P.O. Box 9203, FI-02015 TKK, Espoo, Finland
Harri.Valpola@tkk.fi
<http://www.cis.hut.fi/projects/dss/>

Abstract. The paper presents an example of exploratory data analysis of climate measurements using a recently developed denoising source separation (DSS) framework. We analysed a combined dataset containing daily measurements of three variables: surface temperature, sea level pressure and precipitation around the globe. Components exhibiting slow temporal behaviour were extracted using DSS with linear denoising. These slow components were further rotated using DSS with nonlinear denoising which implemented a frequency-based separation criterion. The rotated sources give a meaningful representation of the slow climate variability as a combination of trends, interannual oscillations, the annual cycle and slowly changing seasonal variations.

1 Introduction

One of the main goals of climate data analysis is to extract the slow components of the data. The classical method of principal component analysis (PCA) is often used for this purpose. However, PCA can be improved by using a frequency-based separation criterion [2,3]. This is the main idea of the proposed method.

The proposed method is based on the use of a frequency-based separation criterion [2]. The main idea is to extract the slow components of the data by using a frequency-based separation criterion. The main idea of EOF can be extended to the proposed method.

In order to extract the slow components (ICA) is a useful method. The main idea is to extract the slow components of the data by using a frequency-based separation criterion. The main idea of ICA is to extract the slow components of the data by using a frequency-based separation criterion.

... be fa ica ... (ee, e.g., [4] f ... d c ...). ICA ba ed ... highe-
 ... de ... a ic a d i ... h i e ec bea ... e i i a i ... ca ica ... a i ...
 ech i e ... ch a he Va i a ... h g a ... a i ... [2]. Se e a a e ...
 a ... ICA i c i a e e ea ch ha e a ead bee ... ade [5,6].

I h i a e , e a a e ea he ... e e ... i g a ... e e e i ... f
 ICA ca ed de ... i g ... ce e a a i ... [7]. DSS i age e a e a a i f a e ...
 h i ch d e ... e ce a i e ... i h e i d e e d e ce a ... i b ... a he ...
 f ... h i d d e c ... e ... h i ch h a e i e e i g ... e i e . The i e e i g e ...
 f h e ... e i e i c ... e d b ... e a ... f a d e ... i g ... c e d ... e . F ... e a ... e ,
 i [8], h e ... ce ... i h ... e i e i e a ... a c i a i ... e e i d e i e d
 ... i g DSS i h i e a ... e i g a d e ... i g . The e a d i g c ... e ... e e ce a
 ... e a e d ... h e e - ... E N i ... S ... h e ... O c i a i ... (ENSO) h e ... e ...
 a d e e a ... h e i e e i g c ... e ... e e e ... a c e d a e ...

I h e ... e e ... , e e DSS i h i e a d e ... i g a h e ... , e ...
 ce i g e ... f c i a e d a a a a i . A i d e f e e c b a d i h e d e ... i g
 ... e i ... e d ... i d e i f h e ... b a c e f h e c i a e ... e . The f ... d
 ... c ... e ... a e f ... h e ... a e d ... i g a i e a i e DSS ... c e d ... e b a e d ...
 ... i e a d e ... i g . The ... a i ... i d e ... ch h a h e e ... a c e d c ... e ... e ...
 ... d h a e d i c ... e ... e c , a .

The e ... a c e d c ... e ... e ... e d ... e ... e e h e b a c e f h e ...
 c i a e h e ... e a a a i e a c ... b i a i ... f ... e d , d e c a d a - i e a ... a ... c i -
 a i ... , h e a ... a c c e a d h e h e ... e a i h d i i c ... e c a c ... e
 U i g h i a ... a c h , h e ... c i a e h e ... e a a e i d e i e d a c e a i b -
 ... a c e f h e c i a e ... e a d ... e h e i e e i g h e ... e a h i d d e i
 h e e a h e ... e a ... e e ... a e f ... d .

2 DSS Method

DSS i a g e e a a g i h i c f a e ... h i ch ca be ... e d f ... d i c ... e i g i -
 ... e e i g h e ... e a h i d d e i ... i a i a e d a a [7]. S i i a ... PCA, ICA ...
 ... h e ... a i ... ech i e , DSS i ba e d ... h e i e a ... i i g ... d e . The ba ic
 a ... i ... i h a h e e a e ... e h i d d e c ... e ... $s(t)$ (a ... ca e d ... ce
 ... fac ...) h i ch a e e e c e d i h e e a ... e e ... $x(t)$ h ... g h a i e a ... a -
 ... i g : $x(t) = \mathbf{A}s(t)$. The ... a i g \mathbf{A} i ca e d h e ... i i g ... a i i h e ICA
 e ... i ... g ... h e ... a d i g ... a i i h e c ... e ... f PCA.

The g a f h e a a ... i i ... e i a e h e ... e ... c ... e ... $s(t)$ a d h e
 c ... e ... d i g ... a d i g ... e c ... (h e c ... e ... f \mathbf{A}) f ... h e b e e d d a $x(t)$.
 I h e c i a e d a a a a i , h e c ... e ... e ... a ... c ... e ... e d ... h e i e -
 ... a i g a e f h e c i a e ... e a d h e ... a d i g ... e c ... a e h e ... a i a ... a ...
 ... h i g h e ... i c a ... e a h e ... a e ... c ... e ... d i g ... h e f ... d c ... e ... e
 The c ... e ... $s(t)$ a e ... a ... a i e d ... i i a i a c e , a d h e e f ... e h e
 ... a i a ... a e ... h a e a ... e a i g f ... c a e .

The ... e ... e ... f DSS i ... -ca e d h i e i g ... h e i g (ee Fig. 1). The
 g a f h i e i g i ... i f ... h e c ... a i a c e ... c ... e ... f h e d a a i ... ch
 a a ... h a a ... i e a ... e c i ... f h e d a a h a ... i i a i a c e . The ... i i e

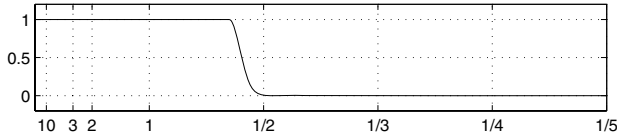


Fig. 2. Frequency response of the filter used in DSS with linear denoising. The abscissa is linear in frequency but is labeled in terms of periods, in years.

... e. f. e.) a 1... d ced 1 [8] ... ide if he. b. ace 1 h. ... 1 e. 1. e. a. ... a. a. lab 1 1 .

Afe. ha, he... c... e... a. e... a. ed. ch ha. he... d ha. e di - 1 c... e... ec. a. Thi 1 d. e b. DSS 1 h. he ge. e a 1 e a 1 e... ced. e 1. hch de. 1 1 g. 1 e e... a. f e... e c -ba. ed. e a a 1. c. 1 e. 1. . P. ac. ica , he de. 1 1 g... ced. e 1 ba. ed... h. e 1 g... . hed di. c. e e c. 1 e . a. - f... (DCT) ... e... ec. a. f he c... e... a. d. 1 g. 1 e. e DCT ... ca. c. a. e he de. 1 ed... ce $f(S)$. Thi de. 1 1 g. echa. 1. 1... e ha. 1 1 a... he h. e 1 g-ba. ed e 1 a 1. f he... ce. a. ia. ce... ed 1 [9]. The a. g. 1 h a... 1 e... de. he... ce. acc. di. g... he. f e... e. c. 1 e 1 g... g. a. h. ic. 1 dea... e ha. 1 1 a... [10].

3 Data and Preprocessing Method

The ... ed ech 1 e 1 a 1 ed ... ea. e e... f hee. a... a... he. ic. a. a. lab e... face e... e a... e, ea. e e... e... e a d... ec. 1 a 1. . Thi... e f a. a. lab e 1. f e... ed f... de. c. 1 b. 1 g. g. ba. c. 1 a. e. he... e a... ch a ENSO [11]. The da. a. e... a. e... ided b... he. ea. a 1. ... ec. f he Na. 1. a. Ce... e. f. E... 1... e a P. ed. ic. 1. Na. 1. a. Ce... e. f. A... he. ic Re. ea. ch (NCEP/NCAR) [12].¹

The da. a. e... e. g. ba. g. ided. ea. e e... e a... g. e. 1 d. f 1 e. The... a. ia. g. id 1. eg. a... aced... e. he. g. be 1 h $2.5^\circ \times 2.5^\circ$ e... 1. . A h. gh. he... a 1. f he da. a 1... e f... he be. g. 1 g. f he... ea. a 1. e. 1 da. d 1 c... ide. ab... a. re. h... gh... he. g. be, e... ed he. h. e. e. 1 d. f 1948-2004.

The... g- e... ea... a... e... ed f... he da. a a d. he da. a 1... e e. eighed 1 1 a... [8] ... di. 1 1 h. he e. ec. f a. de. e... a... 1 g. g. id a... d he... e. Each da. a 1. ... a... 1 ed b... a. eigh... 1. a... he... a. e... f he c... e... di. g. a. ea. f 1... ca. 1. . The... a. ia. di. e. 1. a 1. f he da. a a... ed ced 1 g. he PCA/EOF a... a 1 a... ed... he... eighed da. a. F... each da. a e... e. e. a. 1 ed 100... 1. c. 1 a. c... e... h. ch e... a... e ha. 90% f he... a... a. ia. ce. The DSS a... a 1. a... he... a... 1 ed... he. c... b. 1 ed da. a. c... a 1 g. he... ea... e... e... f he h. ee. a. ia. b. e .

¹ The authors would like to thank the NOAA-CIRES Climate Diagnostics Center, Boulder, Colorado, USA, for providing NCEP Reanalysis data from their Web site at <http://www.cdc.noaa.gov>.

4 Results

Figure 1 identifies the basic frequency characteristics of the DSS 1 high-frequency signal. The 1-yr DSS 1 high-frequency signal is characterized by high-frequency variability. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 3. The analysis of the DSS 1 high-frequency signal is shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4.

Here, the frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 3. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4.

The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 3. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4.

The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 3. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4.

Analysis of the DSS 1 high-frequency signal, the frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 3. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4.

Characteristics of the DSS 1 high-frequency signal are shown in Figure 3. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4. The frequency characteristics of the DSS 1 high-frequency signal are shown in Figure 4.

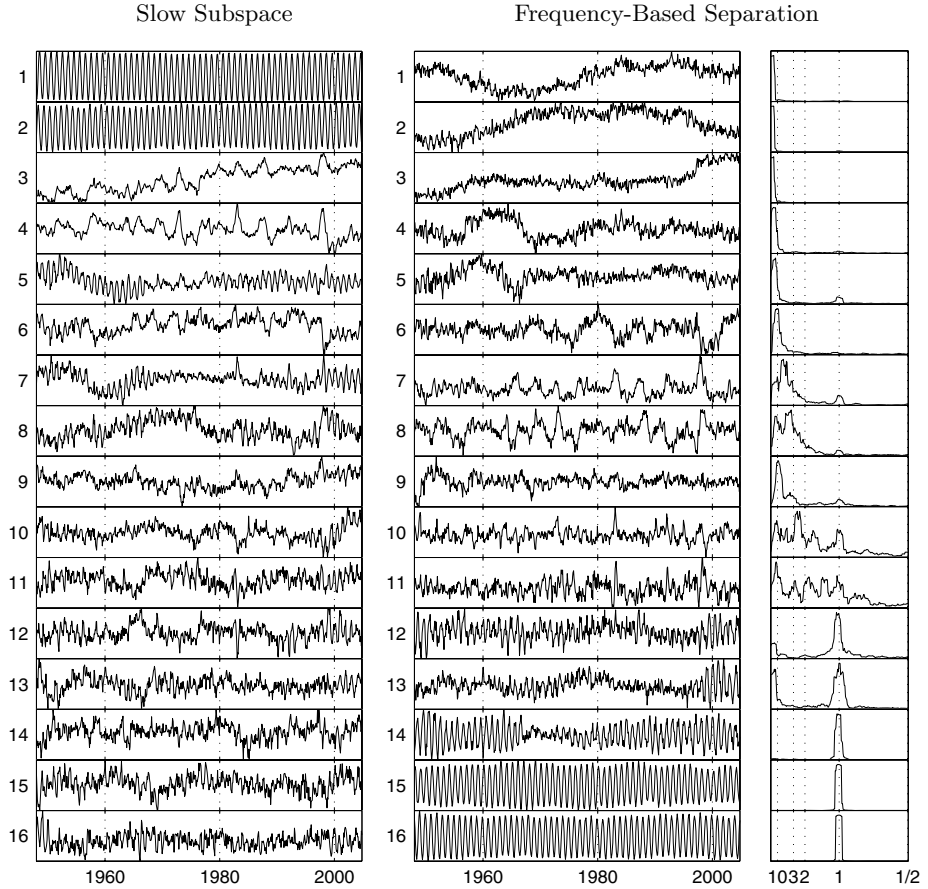


Fig. 3. Left: The monthly averages of the components extracted by DSS with linear denoising. Middle: The rotated slow components estimated by frequency-based DSS. The variances of all the components are normalised to unity. Right: The power spectra of the components found by frequency-based DSS. The abscissa is linear in frequency but is labeled in terms of periods, in years.

1 the N, he, He 1 he, e1 he e ea e a A 1 1a ce a a e ac ed 1 [8].

C...e... 12 16 ha e... 1 e c...e...-a... a fe e cie 1 he1... e... ec a. The a... a c ce... a ea... 1 c...e... 15 16. The e f he... ce ee be he a... a ci a 1... d a ed (... i 1 ed) b... e... ig a. Th... , hi e f c...e... a be e a ed... e he... e a... cha gi g he a... a c ce. H... e e, a... e a e ad... 1 ed... , he f... d... a... 1... 1 hi... hi... b... a... be... e a 1 gf.

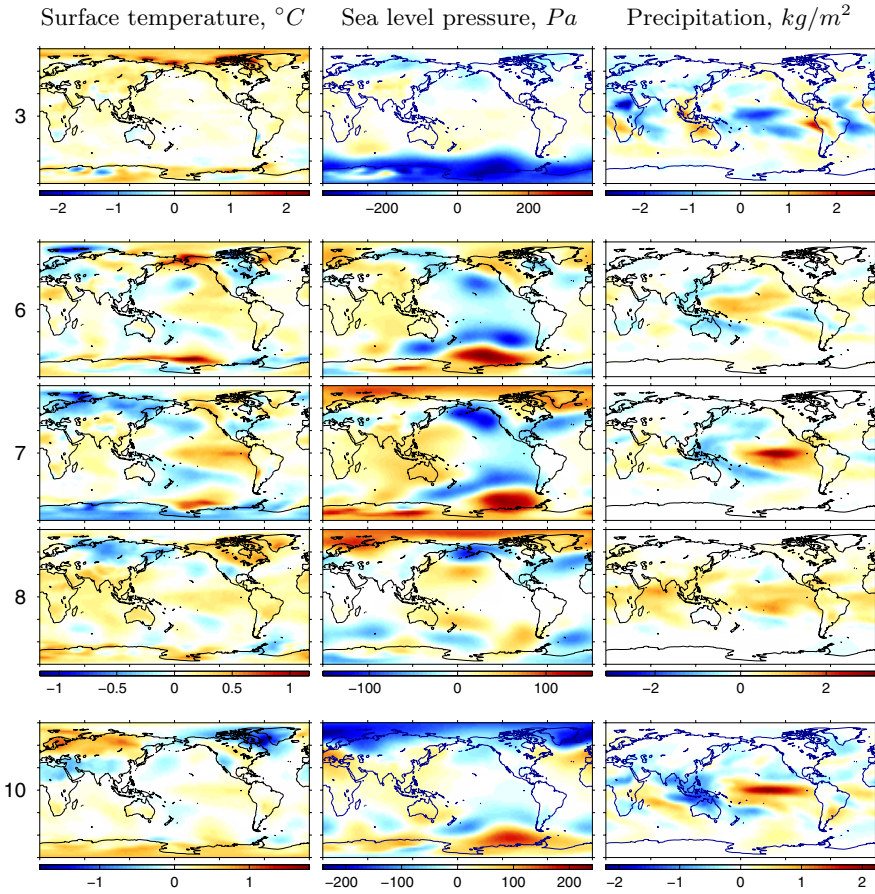


Fig. 4. The spatial patterns of several components found by frequency-based DSS. The label on the left indicates the number of the component in Fig. 3.

5 Discussion

In this paper, we have shown that the DSS method can be applied to a wide range of climate data. We have found several components that are independent of each other and can be used to describe the behavior of the climate system. The results suggest that the DSS method can be used to identify the most important components of the climate system. This is a significant step towards understanding the complex interactions between different parts of the climate system.

Results suggest that the DSS method can be used to identify the most important components of the climate system. This is a significant step towards understanding the complex interactions between different parts of the climate system. The results suggest that the DSS method can be used to identify the most important components of the climate system.

ed. The idea is to separate each of the features of the signal into a set of independent components. The first step is to compute the DSS of the signal. The second step is to separate the signal into a set of independent components. The third step is to separate the signal into a set of independent components. The fourth step is to separate the signal into a set of independent components. The fifth step is to separate the signal into a set of independent components. The sixth step is to separate the signal into a set of independent components. The seventh step is to separate the signal into a set of independent components. The eighth step is to separate the signal into a set of independent components. The ninth step is to separate the signal into a set of independent components. The tenth step is to separate the signal into a set of independent components.

Next, we have to separate the signal into a set of independent components. The first step is to separate the signal into a set of independent components. The second step is to separate the signal into a set of independent components. The third step is to separate the signal into a set of independent components. The fourth step is to separate the signal into a set of independent components. The fifth step is to separate the signal into a set of independent components. The sixth step is to separate the signal into a set of independent components. The seventh step is to separate the signal into a set of independent components. The eighth step is to separate the signal into a set of independent components. The ninth step is to separate the signal into a set of independent components. The tenth step is to separate the signal into a set of independent components.

References

1. H. von Storch and W. Zwiers, Statistical Analysis in Climate Research. Cambridge, U.K.: Cambridge Univ. Press, 1999.
2. M. B. Richman, Rotation of principal components, J. of Climatology, vol. 6, pp. 293-335, 1986.
3. K.-Y. Kim and Q. Wu, A comparison study of EOF techniques: Analysis of non-stationary data with periodic statistics, J. of Climate, vol. 12, pp. 185-199, 1999.
4. A. Hyvärinen, J. Karhunen, and E. Oja, Independent Component Analysis. J. Wiley, 2001.
5. F. Aires, A. Chédin, and J.-P. Nadal, Independent component analysis of multivariate time series: Application to the tropical SST variability, J. of Geophysical Research, vol. 105, pp. 17, 437-17, 455, July 2000.
6. A. Lotsch, M. A. Friedl, and J. Pinzón, Spatio-temporal deconvolution of NDVI image sequences using independent component analysis, IEEE Trans. on Geoscience and Remote Sensing, vol. 41, pp. 2938-2942, December 2003.
7. J. Särelä and H. Valpola, Denoising source separation, Journal of Machine Learning Research, vol. 6, pp. 233-272, 2005.
8. A. Ilin, H. Valpola, and E. Oja, Semiblind source separation of climate data detects El Niño as the component with the highest interannual variability, in Proc. of Int. Joint Conf. on Neural Networks (IJCNN2005), (Montreal, Quebec, Canada), 2005. Accepted.
9. H. Valpola and J. Särelä, Accurate, fast and stable denoising source separation algorithms, in Proc. of Fifth Int. Conf. on Independent Component Analysis and Blind Signal Separation (ICA 2004) (C. G. Puntonet and A. Prieto, eds.), vol. 3195 of Lecture Notes in Computer Science, (Granada, Spain), pp. 65-72, Springer-Verlag, Berlin, 2004.
10. A. Hyvärinen, P. Hoyer, and M. Inki, Topographic independent component analysis, Neural Computation, vol. 13, no. 7, pp. 1525-1558, 2001.
11. K. E. Trenberth and J. M. Caron, The Southern Oscillation revisited: Sea level pressures, surface temperatures, and precipitation, Journal of Climate, vol. 13, pp. 4358-4365, December 2000.
12. E. Kalnay and Coauthors, The NCEP/NCAR 40-year reanalysis project, Bulletin of the American Meteorological Society, vol. 77, pp. 437-471, 1996.

Efficient Processing of Ranked Queries with Sweeping Selection*

Wei Ji¹, Maria Ester¹, and Jianer Han²

¹ School of Computing Science, Simon Fraser University,
{wj, ester}@cs.sfu.ca

² Department of Computer Science, Univ. of Illinois at Urbana-Champaign
hanj@cs.uiuc.edu

Abstract. Existing methods for top- k ranked query employ techniques including sorting, updating thresholds and materializing views. In this paper, we propose two novel index-based techniques for top- k ranked query: (1) indexing the layered skyline, and (2) indexing microclusters of objects into a grid structure. We also develop efficient algorithms for ranked query by locating the answer points during the sweeping of the line/hyperplane of the score function over the indexed objects. Both methods can be easily plugged into typical multi-dimensional database indexes. The comprehensive experiments not only demonstrate that our methods outperform the existing ones, but also illustrate that the application of data mining technique (microclustering) is a useful and effective solution for database query processing.

1 Introduction

Ranked query processing is a fundamental database problem. The answer to a top- k query is the set of k objects with the highest aggregate scores. The combined score function is a real-valued function over the objects. For example, given a database with a score function $f(x, y)$ (where x is the beach and y is the hotel), and the score function $f(x, y) = 0.3x + 0.7y$, the top-3 hotels are the best hotels in the area.

The main challenge in handling top- k queries is to efficiently locate the top- k objects. The main challenge is to efficiently locate the top- k objects. This is a challenging problem because of the large number of objects. Several methods have been proposed, including the use of microclustering [16], the use of the layered skyline [4,8], and the use of the layered skyline [11]. In this paper, we propose two novel index-based techniques for top- k ranked query. The first technique is based on the layered skyline, and the second technique is based on the microclustering. We develop efficient algorithms for ranked query by locating the answer points during the sweeping of the line/hyperplane of the score function over the indexed objects. Both methods can be easily plugged into typical multi-dimensional database indexes. For example, the

* The work was supported in part by Canada NSERC and U.S. NSF IIS-02-09199.

are defined as follows. Let X be a d -dimensional data set, and let $f(x) = \sum_{i=1}^d a_i \cdot x_i$ be a linear function over X , where $a_i \in [0, 1]$ and $\sum_{i=1}^d a_i = 1$. We call f a k -additive function if $f(x) \leq k \cdot \max_{i \in [1, d]} a_i$. A k -additive function f is called a k -additive linear function if $f(x) = \sum_{i=1}^d a_i \cdot x_i$ for some $a_i \in [0, 1]$ and $\sum_{i=1}^d a_i = 1$. We call a k -additive linear function f a k -additive linear function if $f(x) = \sum_{i=1}^d a_i \cdot x_i$ for some $a_i \in [0, 1]$ and $\sum_{i=1}^d a_i = 1$.

The following definitions are taken from [19]. Let X be a d -dimensional data set, and let f be a k -additive linear function. We call f a k -additive linear function if $f(x) = \sum_{i=1}^d a_i \cdot x_i$ for some $a_i \in [0, 1]$ and $\sum_{i=1}^d a_i = 1$. We call a k -additive linear function f a k -additive linear function if $f(x) = \sum_{i=1}^d a_i \cdot x_i$ for some $a_i \in [0, 1]$ and $\sum_{i=1}^d a_i = 1$.

2 Foundations

Let X be a d -dimensional data set, and let f be a k -additive linear function. We call f a k -additive linear function if $f(x) = \sum_{i=1}^d a_i \cdot x_i$ for some $a_i \in [0, 1]$ and $\sum_{i=1}^d a_i = 1$. We call a k -additive linear function f a k -additive linear function if $f(x) = \sum_{i=1}^d a_i \cdot x_i$ for some $a_i \in [0, 1]$ and $\sum_{i=1}^d a_i = 1$.

Lemma 1. Let X be a d -dimensional data set, and let f be a k -additive linear function. If $p \succ q$, then $f(p) < f(q)$.

Theorem 1. Let $K \geq k$. Let L_1, \dots, L_K be k -additive linear functions. Then $f(p) < f(q)$ if $p \succ q$.

Definition 1. (MicroCluster[19]) Let C be a d -dimensional data set, and let n, r be positive integers. We define $CF1(C), CF2(C), CF3(C)$ as follows: $CF1(C) = \frac{CF1(C)}{n}$, $CF2(C) = \frac{CF2(C)}{n}$, and $CF3(C) = \frac{CF1(C)}{n} - (\frac{CF2(C)}{n})^2$.

The following definitions are taken from [19]. Let X be a d -dimensional data set, and let f be a k -additive linear function. We call f a k -additive linear function if $f(x) = \sum_{i=1}^d a_i \cdot x_i$ for some $a_i \in [0, 1]$ and $\sum_{i=1}^d a_i = 1$.

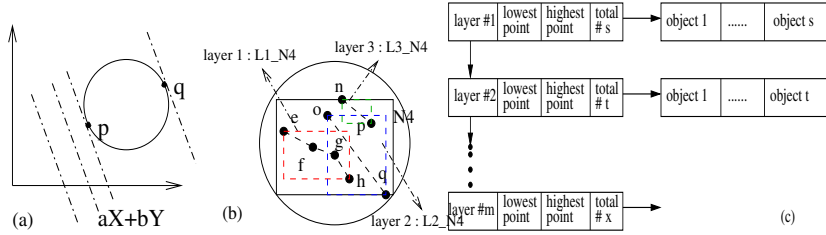


Fig. 1. (a) Contact Points (b) Layered-Skyline in N_4 (c) Linked List for Layered Skyline in a Block

...ic, c...e, 1...eaf...de1...CF-...ee [19]. We c...e the $K(k \leq K)$ a...e...f...1...e...b...c...b...e...c...1...e...a...1...g...e...1...g...1...e...c...1...g...a...g...1...h...1...c...h...a... [14], a...d...1...e...e...e...d...1...e...a...1...1...1...e...e...f...1...d...e...f...1...d...e...c...e...

3 A KNN-Based Sweeping Approach for Top- k Queries

Here we describe the hybrid algorithm for the top- k query. (1) Sweeping over blocks. The algorithm has the I/O complexity $O(n \log n)$ and (2) Sweeping over blocks. The algorithm has the CPU complexity $O(n \log n)$.

(1) Sweeping over blocks. During the sweeping process, the horizontal axis is the block index, the vertical axis is the distance, and the sweep line moves from left to right. Based on the sorted K -nearest neighbors, the distance between a block and its k -th neighbor is the k -th neighbor distance. Each block has a low/high neighbor distance. The algorithm maintains a priority queue Q of the current k -th neighbor distance of the blocks. The algorithm maintains the current distance of the block. The current distance of the block is the current distance of the block. The current distance of the block is the current distance of the block. If the current distance of the block is less than the current distance of the block, the current distance of the block is updated. If the current distance of the block is greater than the current distance of the block, the current distance of the block is updated. The algorithm maintains the current distance of the block.

Algorithm 1 Block-and-Block Ranking (BBR) Method.

Input: k , a data set S in a 2D space.

Output: Top- k nearest neighbors Q .

Method:

1. $Q := R \cup B$;
2. WHILE $|Q| < k$ DO
3. $F := \text{find } k\text{-th neighbor of } Q$;
4. $S := S \cup I(B, F)$; // S is a set of blocks and I is a function.
5. FOR each block $b \in S$ DO

6. IF ... the ha k bec 1 Q ha 1 g... a e, c, e ha s
7. Di ca, d s;
8. ELSE Insert s to Q;
9. O ... k bec f... Q;

Ag 1 h 1 ha e di e e 1 e e a 1 f **SweepIntoBlock**. I ca 1 e ad he b c a d acce a he be g e d bec (e e a **BBR1**). F a MBR 1 R- ee, he e /high e 1 a e he e / e e, gh c, e 1, hie f, a ic, c e 1 CF- ee, he e a e c c ac 1 f he ee 1 gh e, a e he he e, h a p a d q 1 Fig. 1(a). Gi e a ee 1 gh e a e y = $\sum_{i=1}^d a_i \cdot x_i$ a d a ic, c e F 1 h ad i R f he 1 g 1 : $F(x_1, \dots, x_d) = \sum_{i=1}^d x_i^2 - R^2 = 0$ (1). T b a i p a d q, e e e e e e e $\nabla F(x'_1, \dots, x'_d) = c \cdot \mathbf{A}$ (2) ge he 1 h (1). He e c 1 a f ee a i a b e, $\mathbf{A} = a_1, \dots, a_d$, a d $\nabla F(x_1, \dots, x_d)$ h i c h e a he g a d i e f F a X(x'_1, \dots, x'_d), 1 ($\frac{\partial F}{\partial x_1}(x'_1, \dots, x'_d), \dots, \frac{\partial F}{\partial x_d}(x'_1, \dots, x'_d)$).

(2) **Sweeping within layered blocks.** I d e e a a i d c e 1 g e e a b e c 1 e a c h b c, e a e e f he a e e d 1 e 1 c e he g r e a c e e f he d a d i b 1, a d d e e a e c i e e e 1 g 1 h i - a e e d - b c e h d (e e a **BBR2**), a a c e d e SweepIntoBlock 1 Ag 1 h 1. S e e he e 1 a e a f d e h a m a e f 1 e e, he e /high e b e c a e a he a b e f b e c f, h a a e 1 a i a e d. A h e he e d e 1 Fig. 1(b), b e c e, f, g a d h a e a e - 1 e b e c 1 e d e N_4, e a 1 e b e c 1 a e 1 1 1 a h e e c a g e - b e d b a e d - d e d e a L_1 - N_4. The 1 e d 1 e e a g e c e f, he a e e d 1 e 1 h 1 Fig. 1(c), he e he h e d e 1 he e a i a 1 f a 1 f he e d - d e h i c h 1 e 1 b d i g 1 e b e c . N i f a e a f 1 i c h e f he e e, e e e a d he e d - d e h a h e b e e 1 a c c d i g he c e f c i l .

4 A Grid-Based Sweeping Approach for Top-k Queries

A h gh he KNN-ba e d a - a c h c a e c i e b a i he - k b e c , 1 a 1 1 1 a d c a e a he b e c 1 a b c e e he he a e 1 g e c h 1 e 1 a e d. I h i e c 1, e e e a a e a 1 e g i d - b a e d e h d f e e f c i e g a 1 1 g he b e c . S i c e he e - e c i e d e i g h f a c e f c i l 1 f e h a e a f a h e h a a c 1 e a i c, a 1 a e, e e c e 1 g e e b e a c c e a b e f h i a 1 g 1 c a 1 e e d e e 1 e. The b a i c i d e a 1 b i d a g i d - 1 e a 1 1 f he a d a c c e h e b e c 1 h i a b c a g h e g i d. F e CF- ee, a he - g i d a 1 1 a d e (h e R- ee c a e c a b e a d a e d b b d i g a

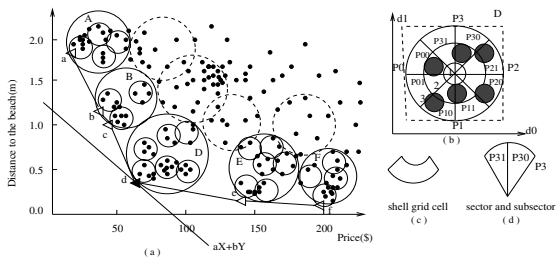


Fig. 2. Shell-Grid Partition of Microclusters

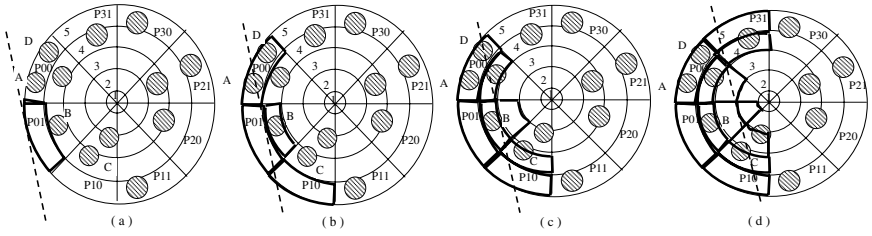


Fig. 3. Sweeping a Shell Grid

MBR... (e.g., $P_{10}, P_{11}, P_{20}, P_{21}, P_{30}, P_{31}$). The bounding rectangle of the grid cell $C_{i,j}$ and the rectangle $R_{i,j}$ are aligned with the grid cell $C_{i,j}$ and the rectangle $R_{i,j}$ is a grid cell $C_{i,j}$ are accepted before the next grid cell $C_{i,j}$. This approach can be used for a set of k rectangles R_1, \dots, R_k and a set of k rectangles R_1, \dots, R_k and a set of k rectangles R_1, \dots, R_k . The bounding rectangle of CF -cell $C_{i,j}$ is the grid cell $C_{i,j}$ and the rectangle $R_{i,j}$ is a grid cell $C_{i,j}$. We call a cell $C_{i,j}$ a shell grid cell (shell cell) if $C_{i,j}$ is a grid cell $C_{i,j}$. We can use the algorithm of the sweep selection to find the bounding rectangle of the CF -cell [13]. More details are given in [1], where the algorithm is described in detail. The algorithm is based on the idea of the sweep selection. The algorithm is based on the idea of the sweep selection. The algorithm is based on the idea of the sweep selection. We will use the notation $2^{d-1} \cdot 2d$ faces of the $(d-1)$ -dimensional hypercube, each having the volume $1/(2^{d-1} \cdot 2d)$ of the $(d-1)$ -dimensional hypercube, face of the cube. The $2^{d-1} \cdot 2d$ sectors P_0, \dots, P_{2^d-1} are defined by the $(d-1)$ -dimensional hypercube, face of the cube (see Fig. 2(b) and Fig. 2(d)). Each P_i is divided into 2^{d-1} subsectors $P_{i0}, \dots, P_{i(2^d-1)}$, and P_{31} and P_{30} (Fig. 2(d)). The hypercube is divided into 2^d shells and 2^d sectors. We have the following lemma:

Lemma 2. Let $x_i, x_j, \dots, x_{i-d}, \dots, x_{j-d}$ be the coordinates of the points $P_i, P_j, \dots, P_{i-d}, \dots, P_{j-d}$. If $0 \leq j < d, i \neq j, |o_i - x_i| \leq |o_j - x_j|$, then $|o_i - x_i| \leq |o_j - x_j|$. If $i \geq d, j \geq d, j \neq (i-d), |o_{i-d} - x_{i-d}| \geq |o_j - x_j|$.

We first consider the case $0 \leq j < d$. Let s_0, \dots, s_{d-1} be the bits of the binary representation of i . If $s_{i-d} = 1$, then $s_{i-d} = 1$ and $s_{i-d} = 1$. For each s_j , if $s_j = 1$, then $x_j > x_i$, and if $s_j = 0$, then $x_j < x_i$. Each s_j has $2^{d-1} - 1$ neighbors. A $(d-1)$ -dimensional hypercube has $(d-1) \cdot 2^{d-2} - 1$ neighbors. A $(d-1)$ -dimensional hypercube has $(d-1) \cdot 2^{d-2} - 1$ neighbors. A $(d-1)$ -dimensional hypercube has $(d-1) \cdot 2^{d-2} - 1$ neighbors. Geometrically, each P_{mn} has a bounding rectangle R_{mn} and a bounding rectangle R_{mn} . A $(d-1)$ -dimensional hypercube has $(d-1) \cdot 2^{d-2} - 1$ neighbors.

Q is defined as the edge e in the grid G such that e is adjacent to h and h is adjacent to e . The set of all such edges is denoted by E_h .

Algorithm 2 A Shortest Grid Ray (SGR) Method

Input: CF, the grid G , and S (a point), k .

Output: The k nearest points T .

Method:

1. Calculate the distance $d(e, S)$ for each edge $e \in E_h$; $Q = \emptyset$; $T = \emptyset$;
2. Insert Q into a priority queue CF ;
3. WHILE CF is not empty and $|T| < k$ DO
4. Remove the edge e from CF ;
5. IF $e \in E_h$
6. insert e into Q and delete its neighbors from CF ;
7. ELSE IF $e \in E_h$ and e is adjacent to h
8. insert e into Q and delete e from CF ;
9. ELSE IF $e \in E_h$ and e is adjacent to h
10. insert e into Q and delete its neighbors from CF ;
11. ELSE IF $e \in E_h$
12. add e to T ;
13. Omit k nearest points T ;

The algorithm is based on the idea of the edge distance $d(e, S)$ where e is the edge in the grid G such that e is adjacent to h and h is adjacent to e (i.e., MC2) and h is the edge in the grid G such that h is adjacent to q and q is adjacent to h (i.e., MC1). The edge e is called a q -edge if $d(e, S) = d(q, S)$. The edge e is called a w -edge if $d(e, S) = d(w, S)$. The edge e is called a p -edge if $d(e, S) = d(p, S)$. The edge e is called a p -edge if $d(e, S) = d(p, S)$. The edge e is called a p -edge if $d(e, S) = d(p, S)$. The edge e is called a p -edge if $d(e, S) = d(p, S)$.

5 Experiments

In this section, we evaluate the performance of the proposed algorithm. The experiments were conducted on a Pentium III 800MHZ machine with 512M RAM, running Windows XP. We used the following parameters

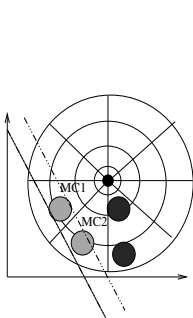


Fig. 4. Error Bound

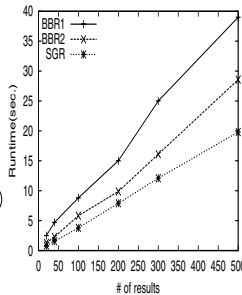


Fig. 5. Query time (1)

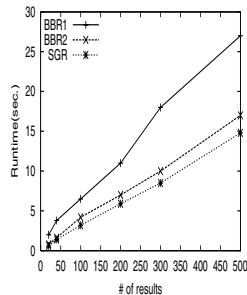


Fig. 6. Query time (2)

and O₁ in C++, and based on PREFER in .db.c.d.ed / .efe. The test set da a e . . f 100,000 records with 5 attributes in the domain /c . . e a e d d i . b i . . e e g e e a e d b h e d a a g e e a . . . f [2].

(1) Comparison of BBR and SGR. Fig. 5 and 6 highlight the algorithm's performance in terms of execution time and coverage rate. BBR2 consistently outperforms BBR1 in terms of execution time, while SGR ($\epsilon = 10$) is the fastest in terms of coverage rate. When ϵ changes from 5 to 10 and 15, the execution time of da a e i e d e e d e h e d e e a i g b e f f i i e d i c c e e a d h e e a e a e a c e a g e a e b e c e e a i e h i g h e d e h e i c e a i g i e f i c c e e . (Fig. 7, 8 and 9).

(2) Comparison with Onion and PREFER. We compare the performance of K against the algorithms BBR, SGR, and K against the algorithms O₁ ($K = 200$). When the dimensionality is from 2 to 5, SGR is the fastest (Fig. 10). BBR2 consistently outperforms O₁ in terms of execution time, while SGR is the fastest in terms of coverage rate. When k changes, the execution time of BBR2 and SGR is higher than that of PREFER (Fig. 11, 12). When k increases, the coverage rate of the algorithms SGR is higher than that of PREFER (Fig. 13 and 14). Because the execution time of da a e b e f f i e h e i e f f a e i a i e d i e f f P R E F E R i s b e d c e d h i c h e a d h e e c e a g e a e h a h e i d e e d e d a a e

6 Related Work

The $-k$ algorithm is based on the algorithm proposed by Fagin in the context of multidimensional data access [7], and it can be categorized into three types: **(1) Sorted accessing and ranking** algorithm, which is based on the algorithm proposed by each of the algorithms in each attribute in the $-k$ algorithm [17]. The algorithm is based on the algorithm proposed by [7]. Finally, the algorithm (TA) [8] is designed to access the data in the $-k$ algorithm. **(2) Random accessing and ranking** algorithm, which is based on the algorithm proposed by each of the algorithms in each attribute in the $-k$ algorithm [10]. The algorithm is based on the algorithm proposed by each of the algorithms in each attribute in the $-k$ algorithm.

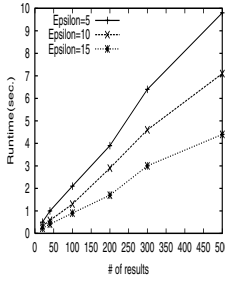


Fig. 7. Effect of Eps(1)

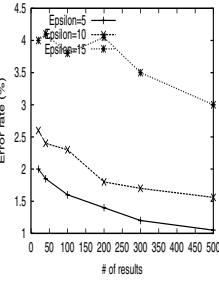


Fig. 8. Effect of Eps(2)

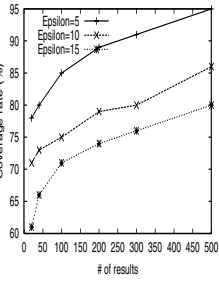


Fig. 9. Effect of Eps(3)

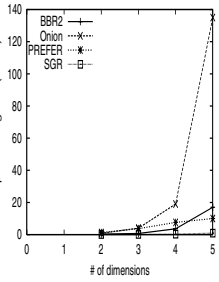


Fig. 10. Preprocessing

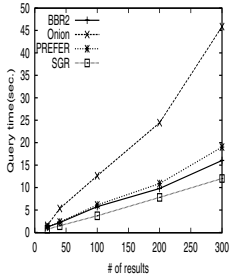


Fig. 11. Query(1)

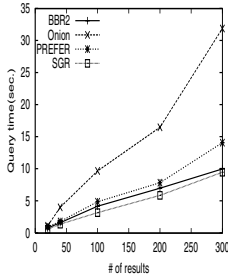


Fig. 12. Query(2)

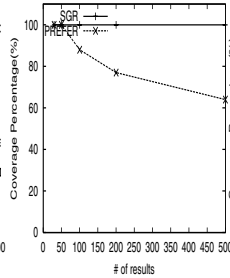


Fig. 13. Quality(1)

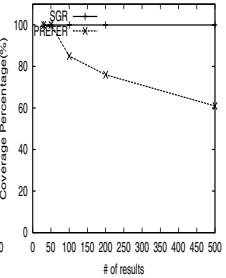


Fig. 14. Quality(2)

... a ... be a he ... c ... e fec ... a chi g ... be , he ea [5] ... e ...
 ... a ... e he ... -k ... e ... i ... a ... a ge ... e ... i ... da aba e. **(3) Pre-materialization
 and rank indices** ... ga i e he ... e i ... a ... ecia a , he a ... ie ... i ... a i ... a ch
 f ... he a ... e ... f ... a ed ... e. I [4], a i ... de i b i ... a ... e ed c ... e h ... f ...
 a ... i ... S ch i de f ... a ge da aba e i e ... e i e d e ... he c ... e h ... d i g
 c ... e i ... [11,12] ... e ... e- a e i a i e i e ... f d i e e ... -k ... e ... a ... e ...
 W he ... he ... e ... c ... e f ... c i ... i ... c ... e ... a ... i e ... ? , a ... a ... be ... f ... he ... e i ...
 h a ... i e ... e ... c a ... f ... he ... e ... e , a d he ... e ... e ... c a ... be ... d ... c e d i ... a
 i ... e i ... e d ... f a h i ... A ... he ... e i ... g ... a ... e e h ... a ... a ... e ... h ... d ... e a ... c h ... e
 f ... e ... , i ... f ... e ... e ... he ... h ... e d a ... a ... e i ... e a c h ... e ... [16] ... e ... e ... a ... e d i ... d e
 ... e ... e ... e ... -k ... e ... (k ≤ K) b ... i ... a ... i e ... e ... d i ... e ... i ... a d c a ... h a ... e a
 h g e ... be ... f ... a ... e i a i e d ... a ... i

7 Conclusions

Ra -a a e ... e ... c ... e i g ha ... e c ... e ... e e g e d a a i ... a ... a a d i g i
 da aba e ... e ... , a d ... f ... e ... i ... g ... e h d e ... i ... a ... e i a i a i ... a d i d e
 ... c ... e. I ... h i ... a ... e ... , e ... e ... e i ... d e i g a e e d ... i ... e a d he -g i d ... i c ...
 c ... e ... f ... -k ... a ... e d ... e , a d ... e ... e ... e h d ... e e i g h e h ... e a ... e f ... he
 ... c ... e f ... c i ... e ... he i ... d e d ... b e c O ... e h d ... c a ... be e a i ... a d a e d ... e -
 i ... i g ... r -d i ... e ... i ... a i ... d e ... c ... e. The e ... e ... i ... e ... a ... e ... d e ... a ... e ... he
 ... e g h f ... e h d a d he ... e f ... e ... e ... f ... he ... i c ... c ... e i g e c h i ... e i ... -k
 e ... e ... c ... e i g .

References

1. S. Berchtold, C. Bhm and H. P. Kriegel. The Pyramid-Technique: Towards Breaking the Curse of Dimensionality. *SIGMOD* 1998.
2. S. Borzsonyi, D. Kossmann, and K. Stocker. The Skyline Operator. *ICDE* 2001.
3. N. Beckmann, H. P. Kriegel, R. Schneider, and B. Seeger. The R*-tree: An Efficient and Robust Access Method for Points and Rectangles *SIGMOD* 1990.
4. Y. C. Chang, L. D. Bergman, V. Castelli, C. S. Li, M. L. Lo and J. R. Smith. Onion Technique: Indexing for Linear Optimization Queries. *SIGMOD* 2000.

5. S. Chaudhuri and L. Gravano. Evaluating Top-k Selection Queries. *VLDB* 1999.
6. T. Cormen, C. E. Leiserson, et al. Introduction to Algorithms, The MIT Press, 2001.
7. R. Fagin. Fuzzy Queries in Multimedia Database Systems. *PODS* 1998.
8. R. Fagin. et al. Optimal Aggregation Algorithms for Middleware. *PODS* 2001.
9. A. Guttman. R-trees: A dynamic index structure for spatial searching. *SIGMOD* 1984.
10. S. Guha, et al. Merging the Results of Approximate Match Operations. *VLDB* 2004.
11. V. Hristidis, N. Koudas, and Y. Papakonstantinou. PREFER: A System for the Efficient Execution of Multi-parametric Ranked Queries. *SIGMOD* 2001.
12. V. Hristidis and Y. Papakonstantinou. Algorithms and applications for answering ranked queries using ranked views. *VLDB J.* 2004.
13. Wen Jin, Jiawei Han, Martin Ester Mining Thick Skylines over Large Databases. *PKDD* 2004.
14. D. Papadias, Y. F. Tao, G. Fu and B. Seeger. An Optimal and Progressive Algorithm for Skyline Queries. *SIGMOD* 2003.
15. N. Roussopoulos, S. Kelley, and F. Vincent. Nearest Neighbor Queries. *SIGMOD* 1995.
16. P. Tsaparas, et al. Ranked Join Indices. *ICDE* 2003.
17. E. L. Wimmers, L. M. Haas, M. T. Roth and C. Braendli. Using Fagin's Algorithm for Merging Ranked Results in Multimedia Middleware. *CoopIS* 1999.
18. D. A. White and R. Jain. Similarity Indexing with the SS-tree. *ICDE* 1996.
19. T. Zhang, R. Ramakrishnan, and M. Livny BIRCH: an efficient data clustering method for large databases. *SIGMOD* 1996.

Feature Extraction from Mass Spectra for Classification of Pathological States

Alejandro Kalousis, José Prados, Emilio Rehner, and Mercedes Hilari

University of Geneva,
Computer Science Department,
Rue General Dufour, 1211, Geneva, Switzerland
{kalousis, prados, hilario}@cui.unige.ch, rexhepaj@unil.ch

Abstract. Mass spectrometry is becoming an important tool in proteomics. The representation of mass spectra is characterized by very high dimensionality and a high level of redundancy. Here we present a feature extraction method for mass spectra that directly models for domain knowledge, reduces the dimensionality and redundancy of the initial representation and controls for the level of granularity of feature extraction by seeking to optimize classification accuracy. A number of experiments are performed which show that the feature extraction preserves the initial discriminatory content of the learning examples.

1 Proteomics and Mass Spectrometry

Clinical and biological changes in protein levels, degradation, and modification are associated with disease. Mass spectrometry is a powerful tool for the identification of these changes. The identification of these changes, however, is a challenging task due to the high-dimensionality and high redundancy of the data.

In this paper, we present a new feature extraction method that aims to reduce the dimensionality and redundancy of the data while preserving the discriminatory content. The method is based on a greedy search for a set of features that maximizes the classification accuracy.

One of the main challenges in the analysis of mass spectrometry data is the identification of the changes in protein levels, degradation, and modification. The identification of these changes is a challenging task due to the high-dimensionality and high redundancy of the data. In this paper, we present a new feature extraction method that aims to reduce the dimensionality and redundancy of the data while preserving the discriminatory content. The method is based on a greedy search for a set of features that maximizes the classification accuracy. [1]. The method is based on a greedy search for a set of features that maximizes the classification accuracy. [2].

The a_{e1} , g_{a1} ed a f... ec 1. 2 de cibe he... ce 1 g ha...
 ea... a... ec... e... be; ec 1. 3 e... hich... f...
 ce 1 g, a d h... a e e a ed i h fea... ac 1. a d h... c... he
 di e... a 1... he... e... e... a 1...; ec 1. 4 e... e... g a e he deg ee
 f di e... a 1... ed c 1. b... gh b fea... ac 1... e e f... a e i e
 f c a 1 c a 1... e... e... i... de... e a b i h... c a 1 c a 1... e f... a ce,
 a d e h b i h... e c a... a i c a... e e c... he a... i a e... a a e... a e
 f... e... ce 1 g a d fea... ac 1...; a... e c... c... de 1... ec 1. 5.

2 Mass Spectra Preprocessing

A... ec... f a b i... g a... a... e... e... d i... e... a... i g a. The -a 1
 c... e... d... he... a... /... a... e... f... e... i... de... e... c... ed... he b i... g a... a... e
 a d he -a 1... he 1... e... i... f... he... e... a... e... he a... e... i... g... g... e... a... ed... he
 c... c... e... a... i... f... he c... e... e... d i... g... e... i... h... e... a... e... A... a... ec... i... a
 e... c... h... e... d i... e... a... i... e... a... he... b... e... f... /... a... e... ec... d... e... d... b... he
 e... c... e... e... he... a... e... f... each... di... e... i... h... e... i... e... i... f... he c... e... e... d i... g
 /... a... e... I... e... i... e... f... e i g h b... i... g... /... a... e... a... e... h i g h... c... e... a... ed... e... i... g
 1... h i g h... a... i... a... e... d... d... a... c... a... i... g... he fea... e... f... a... a... ec... .

M... ec... a... d... a... d... c... i... d... e... a... b... e... e... f... e... e... ce 1 g, h i c h c a n b e
 g h... d i... d... : b a e i e... e... a... , d e... i... i... g... ,... h i g... ,... a... i... a... i... , e a
 d e c 1. a d c a i b... a... i... , [1]. We 1 de cibe... h... e... ac... ed... each... f... he... .

The b a e i e... a... e... f... he 1... e... i... e... f... a... e... , h i c h h a... e... a... i...
 a... e... a... e... , a d... a... i... e... b... e... e... d i... e... e... c... a... I... d... e... f... c... a... i...
 b... e... e... i... e... i... e... f... /... a... e... b... e... e... a... i... g... f... i... h... a... b... e... b... a... c... e... d... T...
 c... e... he b a e i e... e... d... a... c... a... e i g h... e... d... a d... a... c... i... g... ,... he 1... f...
 he... c... a... i... a... e... a... c... e... d... f... he... e... c... . O... he... e... e... d... a... e... f... c... a...
 i... i... a... a... e... e... a... c... h... f... c... a... i... a... a... e... f... e... d... , he... e... i... g... ,... h...
 ... a... a... a... i... . U... i... g... he... e... c... a... i... a... he... i... g... a... i... i... ,... r... e... c... e... i... e
 c... a... a... a... a d... he... a... b a e i e... i... i... c... e... d... b... he... e... a... i... c... a... i... f...
 he... i... a... c... a... e i g h... e... d... a d... a... c... i... g... ,... he... r... e... c... e... c... a... i... g... a... .

T... d... e... i... e... a... d... h... he... i... g... a... e... d... a... e... e... d... e... c... a... i... i... c... e... d
 i... h... a... e... d... i... e... ; a d... e... a... e... d... d... e... c... i... i... i... g... r... e... i... e... c 1. 3. S i g a... i... e... i... e
 a... e... e... a... e... d... ,... b... e... e... d... e... d... e... e... e... i... e... a... c... d... i... ,... i... a... a... i...
 c... e... h i c h... i... e... i... a... e... e... a... i... g... i... h... he L_1 ... f... he... e... c... .

P... a... d... e... c 1. 1... he... d... e... c 1. f... c... a... a... i... a... i... he... a... e... c... . A... e... a
 c... e... c... i... e... e... e... a... he... /... a... e... h... a... d... e... i... ,... h... a... i... :... a... i... g... f... i...
 e... f... c... e... c... a... i... i... a... d... i... g... i... i... g... h... c... e... c... a... i... i... . The
 i... e... i... e... f... a... he... e... i... g h b... i... g... /... a... e... e... h i b i... a... h i g h... e... e... f... e... d... d... a... c... ,
 h... b... e... e... e... i... g... a... e... c... . i... a... i... e... a... e... c... i... d... e... a... b... e... d... c... e... h... e
 e... e... f... a... i... a... e... d... d... a... c... . P... e... a... c... a... i... b... a... i... e... a... b... i... h... e... h i c h... e... a... a... i... g
 d i... e... e... e... c... a... c... e... d... d... he... a... e... e... a... , i. e. he... a... e... e... i... . We... d... e... d
 he... a... e... a... c... h... f... e... d... i... [3] h i c h... i... e... e... i... a... c... e... e... i... a... g... e... h... i... e... a... c... h... i... c... a
 c... e... i... g... i... h... e... a... d... d... i... a... d... a... i... c... a... i... . The... a... c... e... c... a... i...
 a... e... f... d i... e... e... e... c... a... h... a... c... e... d... d... he... a... e... e... a... .

3 Feature Extraction

Feature extraction is the process of identifying the most relevant features from a large set of data. In this paper, we use a set of features to describe the performance of the system. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance.

We use a set of features to describe the performance of the system. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance.

We use a set of features to describe the performance of the system. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance.

The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance.

4 Experimentation

We used a set of features to describe the performance of the system. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance. The features are derived from the system's output and are used to evaluate the system's performance.

The ea... g... h... ed... a... a... e... g... h... J48, ... h... M=2, C=0.25, a... e... e... eighb... g... h... , IBL, a... d... e... e... e... ach... e... g... h... , SMO, ... h... a... e... e... e... e... a... d... C=0.5. The... e... e... a... . . . f... h... e... g... h... . . . e... h... e... f... h... e... WEKA... ach... e... e... g... e... . . . e... , [7]. The... h... e... e... a... g... g... h... . . . e... e... h... e... . . . ha... h... e... e... e... a... d... e... e... e... f... e... a... g... b... i... e... . Pe... f... . . . a... c... e... e... i... a... . . . a... d... e... . . . g... 10-f... d... . . . a... i... e... d... c... . . . -... a... i... d... a... . . . i... g... f... . . . i... g... i... c... a... . . . d... i... e... e... c... e... . . . i... g... McNe... a... . . . e... . . . h... a... e... e... . . . f... i... g... i... c... a... . . . c... e... f... 0.05.

We... e... . . . ed... h... e... e... :... h... e... deg... e... e... f... d... i... e... . . . i... a... . . . ed... c... i... . . . ach... i... e... d... b... h... e... fea... t... e... e... ac... i... . . . e... cha... . . . ,... h... e... a... f... d... i... c... i... i... a... . . . i... f... . . . a... e... e... . . . ed... b... h... e... d... i... e... e... e... e... . . . f... e... e... . . . ce... i... g... . . . a... e... . . . de... . . . i... g... a... d... fea... t... e... e... . . . ac... i... . . . , a... d... . . . a... . . . h... . . . e... . . . ce... i... g... c... . . . d... b... e... . . . i... e... d... .

We... a... . . . ed... h... e... e... h... e... h... d... f... . . . 0.5... 0.95... i... h... a... e... . . . f... 0.05, a... d... f... . . . 0.95... 0.99... i... h... a... e... . . . f... 0.01. The... deg... e... e... f... d... i... e... . . . i... a... . . . ed... c... i... . . . a... ge... f... . . . 60%... 95%... f... h... e... i... i... a... be... f... fea... t... e... i... h... e... c... . . . e... e... a... . . . ec... . . . de... e... d... i... g... . . . h... e... h... e... h... d... a... d... h... e... d... a... e... . . . D... e... . . . ac... f... . . . ace... e... i... e... f... h... e... e... . . . i... . . . a... b... e... 2. The... d... i... e... . . . i... a... . . . ed... c... i... . . . i... d... e... . . . i... ch... a... . . . a... . . . ha... i... e... e... c... d... . . . a... edge... a... d... ed... ce... h... e... . . . a... i... a... . . . ed... da... c... . . . f... h... e... i... i... a... . . . e... e... a...

Table 1. Description of mass spectrometry datasets considered

dataset	# controls	#diseased	mass range (Daltons)	# features
ovarian	91	162	0-20k	15154
prostate	253	69	0-20k	15154
stroke	101	107	0-70k	28664

Table 2. Feature reduction for different values of the wavelet threshold. For each dataset and wavelet threshold, θ , we give: the number of features after feature extraction (# features), and the percentage of feature reduction (reduction %).

θ	prostate		ovarian		stroke	
	# features	reduction %	#features	reduction %	#features	reduction %
0.5	3779	75.06	1591	89.50	11983	58.19
0.6	3538	76.65	1371	90.95	11294	60.59
0.7	3223	78.73	991	93.46	9780	65.88
0.8	2616	82.73	865	94.29	6954	75.73
0.9	1668	88.99	775	94.88	3154	88.99
0.99	1009	93.34	668	95.59	1255	95.62

The... h... e... h... e... a... . . . e... . . . f... e... . . . ce... i... g... e... e... h... e... d... i... c... i... i... a... . . . i... f... . . . a... . . . e... e... a... a... e... d... h... e... ea... i... g... g... h... h... e... d... i... e... e... e... e... a... f... h... e... c... a... i... c... a... be... : 1) h... e...

1) $111a$ c... e... a... ec... a... he... e... ef... edbaeiee... aad...
 ... a1a1... , ... , hi1a1gedaae; 2) hec... ee... a... ec... a... he... e...
 ba... e... e... a... a1a1... , i... e... a... ihdiee... a... ee... h... h... d...
 a... d... hi... ga... e... f... ed, hi1a... g... f... da... e... c... ecie... ide... id... ed
 a... , eachda... e... c... e... d... a... ecic... a... ee... h... h... d; 3) he... da... a... e...
 ... d... cedaf... e... fea... ee... ac1... , c... ecie... ide... id... eda... ,
 ... ide... ih... a... e... f... a... ceba... e... i... ce... i... c... a... a... he... 111a... a... a... abe
 1... f... a... 1... C... a... 1... g... he... e... f... a... ce... f... ea... 1... g... .. a... d... .. e... ca... ee
 h... de... 11... g... a... d... .. hi... g... a... ec... he... di... c... 1... 1... a... .. c... e... .. f... he... ea...
 1... ge... a... e... , h... ec... .. a... 1... .. be... ee... , .. a... d... .. a... .. e... ab1h
 he... e... ec... f... fea... ee... ac1... .. he... di... c... 1... 1... a... .. c... e... .. The... e... 1... a... ed
 e... f... a... ce... (acc... acie...)a... e... g... i... e... 1... g... e... 1.

A... c... ee... a... 1... a... 1... f... g... e... l... h... .. ha... 1... ge... e... a... he... acc... acie... f... he
 ea... 1... g... a... g... 1h... .. he... .. a... d... .. e... ee... a... 1... .. a... e... 1... 1... a... .. he... ba... e...
 1... e... acc... ac... .. . The... e... 1... .. c... ea... .. e... d... a... .. cia... ed... ih... he... e... e... f... de...
 .. 11... g... a... d... .. e... a... ic... di... e... ce... ha... .. d... h... .. a... ce... a... ad... a... age... .. di...
 ad... a... age... f... de... 11... g... .. hi... g... a... d... fea... ee... ac1... ..

The... ab1h... a... ecie... ic... e... f... he... e... ec... f... de... 11... g... .. hi... g... a... d... fea... ee...
 e... ac1... .. di... c... 1... 1... a... .. c... e... .. ec... .. ed... he... 1... g... i... ca... ce... ee... f... he...
 acc... ac... di... e... e... ce... .. he... .. e... ee... a... 1... .. a... d... .. each... f... he... .. , a... d...
 , e... ee... a... 1... , i.e., f... each... a... e... f... he... a... ee... h... e... h... d... , he... e... .. a... e...
 a... 1... ed... 1... abe3... 1... e... .. f... 1... g... 1... ca... .. 1... .. a... d... .. e... .

Table 3. Significant wins and losses table summarized over the different threshold values. A triplet $w/t/l$ for a pair of representations x vs y gives the number of significant wins (w), ties (t), and significant losses for x .

	<i>bl-tic vs all</i>			<i>bl-tic vs peaks</i>		
	SMO	J48	IBL	SMO	J48	IBL
ovarian	0/14/0	0/14/0	1/13/0	0/14/0	8/6/0	0/14/0
prostate	6/8/0	0/14/0	0/14/0	0/14/0	1/11/2	0/14/0
stroke	2/12/0	0/14/0	0/14/0	1/13/0	0/14/0	1/13/0

De... 11... g... a... d... .. hi... g... 1... ge... e... a... .. ee... e... he... di... c... 1... 1... a... 1... g... 1... f... .. a... 1...
 c... .. a... ed... ih... he... ea... 1... ge... a... e... , abe3... , c... H... ee...
 he... e... a... e... h... e... h... d... a... e... f... .. h... ich... ca... 1... ca... 1... acc... ac... 1... g... 1... ca... .. de... e... 1...
 .. a... e... c... .. a... ed... .. he... ba... e... i... e... a... fac... ha... ca... f... .. a... 1... f... .. ed... a... l... f... e... ec... 1... g...
 he... a... .. 1... a... e... h... e... h... d... a... e... A... 1... 1... a... .. ic... .. ea... 1... e... he... ee... ee... a... 1... e... he...
 e... f... a... ce... f... fea... ee... ac1... .. f... he... di... e... e... .. a... e... .. f... he... a... ee... h... e... h... d...
 .. d... , abe3... , c... , ... I... f... he... ca... e... fea... ee... ac1... .. e... a...
 a... di... c... 1... 1... a... .. c... e... .. 1... 1... a... .. ha... f... he... .. e... ee... a... 1... .. ee... he... e...
 he... e... a... .. he... .. g... ch... ice... f... he... a... ee... h... e... h... d... ca... e... ad... .. a... 1... g... 1... ca... .. d...
 1... ca... 1... ca... 1... acc... ac... .. c... .. a... ed... ih... he... ba... e... i... e... acc... ac... ..

The... 1... 1... a... a... e... f... he... a... ee... h... e... h... d... de... ed... he... da... a... e...
 b... .. a... .. he... ea... 1... g... a... g... 1h... .. ed... A... g... d... .. bad... ee... c... 1... f... he... a... ee...

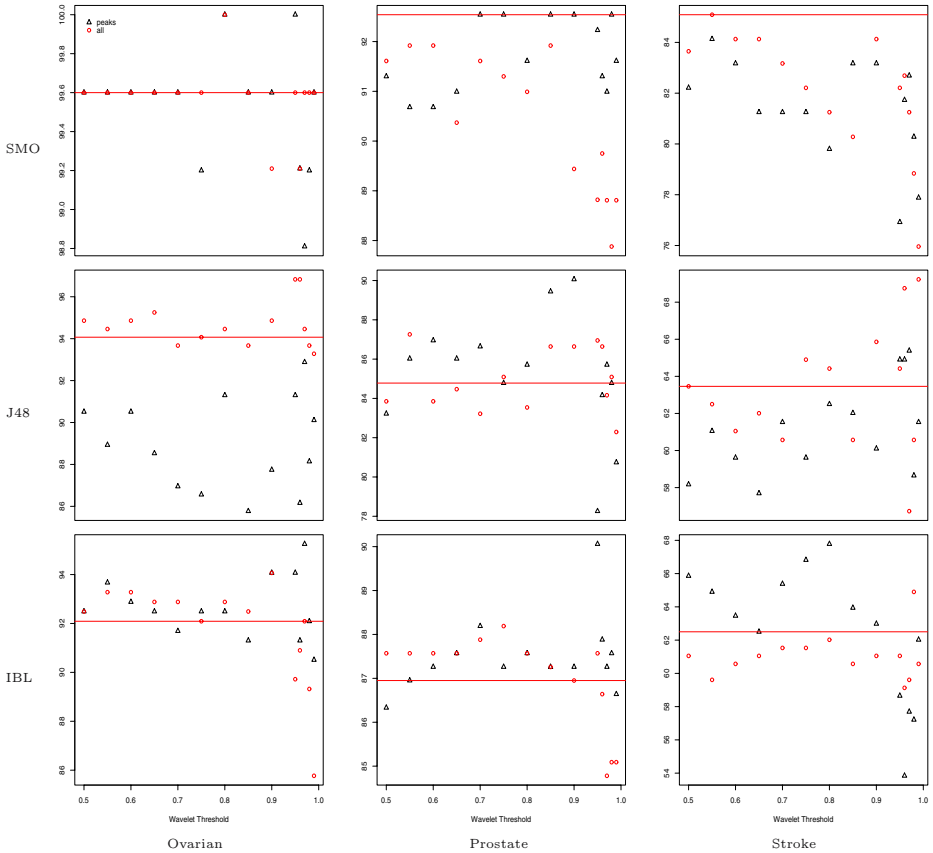


Fig. 1. Accuracy results for each dataset and learning algorithm. Each graph gives the accuracy of a given algorithm for all wavelet thresholds using: the complete spectrum, *all*, only the peaks, *peaks*. The horizontal lines correspond to the accuracy on the *bl-tic* version of the dataset. The y-axis gives the accuracy %, note the differences in scales.

he h d ca ead l a e f a ce ha l e ec l e i g i ca be e e e e ha he ba e l e e f a ce. Th a he e a e a g a l he l e e f he i a e e c l e f he a e e h e h d l e d e l e i l e f e a e e e a c l e. C e a h i e e c l e h d b e g i d e l e l b h e d a a e b a l b h e e a l g a g l h h a e a e a l g l e f c a l l c a l l.

We l d i e e e c h a a e e h e h d h a l d h e h e h i g h e c h a c e l f a l l i g c a l l c a l l e f a c e l h e e i g l a c e f a g i e e a l g a g l h. T a c h e e h a e i g h l e g a e h e h e e e e c e i g l e i e e h e e a l i g l e c e a d a d a c l e l e d e c h l e f a a e e l i g h a l b a e d l c a l l a i d a l l. M e e e c e h e l e e g r e a l l i a l l f l e e i g a a e e a e a d e e c a e a l g a g l h. Th a a a l f

he ,ai 1 g ha e he ,e , ce 1 g a d he ea 1 g a g ,i h a e c , a i -
 da ed (e f d) igh c , ed f , each a e . The a e ha ga e he hige
 acc , ac 1 ch e , he ,ai 1 g e i , e , ce ed i h he ch e , a e , he
 ea 1 g a g ,i h i , ai ed , he ,e , ce ed , ai 1 g e a d e ed , he
 e , e .

E a a i , a d e i h e f d , ai ed c , a i da i . The h e h d
 a e a , g h i c h e e c i , i e f , ed a e he a e a i he ,e i , e -
 e i e . F , each da a e aga i , e , e e a i , e e c , i de ed: , ha
 c , ai a he fea , e he e , de , i i g a d , , h i g i e f , ed , a d
 , , he e fea , e e , ac i , a a , e f , ed , e , a e g i e i a b e 4 .

De , i i g i h a , a i c e e c i , f he a , , i a e a e e h e h d e -
 ai he d i c i i a , i f , a i , f he ea 1 g e a e he c , a e d
 i h he , , e , e e a i . F , a e a e a d a da a e , c a i c a i
 e f , a c e , he , e , e e a i , a e e , i i a , i h h e , he ,
 , e , e e a i , , a i c a i g i c a i d i e e c e a e b e , ed (a i c a
 i g i c a c e da a , h) . The a e h d he e c , a e h e c a i c a i
 e f , a c e f fea , e e , ac i , i h he c a i c a i e f , a c e , , ,
 e f , a c e a e i i a a d , a i c a i g i c a i d i e e c e i b e , ed .

Table 4. Classification accuracy with automatic parameter selection. Automatic parameter selection is only performed for *all* and *peaks*, *bl-tic* is repeated for comparison reasons.

dataset	bl-tic			all			peaks		
	SMO	J48	IBL	SMO	J48	IBL	SMO	J48	IBL
ovarian	99.60	94.07	92.09	99.20	95.25	96.04	99.20	92.49	92.88
prostate	92.54	84.78	86.95	92.23	86.95	86.95	92.54	84.78	90.06
stroke	85.09	63.46	62.50	81.73	65.86	56.73	81.73	62.01	62.01

O e a a , a i c e e c i , f he a e e h e h d a i d he i fa , f
 a a e e c i . I , ed ce he cha ce f e e c i g a h e h d a e ha , d
 e , i a i g i c a i d e e i , a i , f he c a i c a i e f , a c e . E a
 i , , a i e i i a e he e e d f , a i a a d a i a i e i , e c i , f he
 e , e f d e , i i g i , d e , e e c he a , , i a e h e h d , h , e i e i g h e
 a a , f , a i g i c a i b , d e , h i e , he a e i e i e a c e a a i a i e
 a , , a c h (i a i , e c i) i h a b e c i e c i e i , (c a i c a i , a c c , a c) .

5 Conclusion

Ma , e c , e , da a a e cha , a c e i ed b , e , h i g h d i e i , a i a d e ,
 h i g h e e , f , ed , da c a , g h e i fea , e . I h i a e , e e i ed d -
 a i , , e d g e i , d e , e d c e d i e i , a i a d , he a e i e e , e
 , ed da c f , he i i i a , e , e e a i , b fea , e e , ac i . O e f ,
 c e , a , , a a e e d e , , i i . Wa e e , ha e b e e , ed b e f , e i

... a ... ec ... e ... [8,9] a d i ... he d ... a l ... i e ... ea ... i g i ... e ... i e ... i a ...
 11e [10]. Ne e ... he e ... he a ... ed i h ... he a e e c e c i e ... We e ...
 ... he 1 1 1 a ... e ... e a 1 ... f he i g a a d e ... ac fea ... e f ... ha , h ... he
 e ... ac ed fea ... e a e f d i e c b i ... g i c a ... e e a c e . [1] ... e a e e ... e ... ac
 ... ea f ... a ... ec a b ... he d i d ... f ... i h c a 1 1 c a 1 , ... e ... e
 he ... e ed he ... b e ... f he a ... i a e e e f d e 1 1 g a d a a e e
 ... e e c 1 . T a d d e ... he e 1 ... e e i g h ... c ... ed fea ... e ... ac 1 ... i h c a -
 1 c a 1 ... a d ... i d e d a e e c i e a d a ... a i c a ... c ... he g a ... a 1
 ... f fea ... e ... ac 1 ... i h a i e ... a 1 1 1 g c a 1 c a 1 ... acc ... ac .

We h ... d ... e he e ha ... ha ... e ... e 1 ... a fea ... e ... e e c 1 ... e h d .
 O ... g a 1 ... 1 1 1 e he ... b e ... f fea ... e ... ha c a b e ... e d ... e e c 1 e
 ... e f ... c a 1 c a 1 ... b ... e ... ac a h i g h e e , ... e c ... ac , e ... ed d a
 a d e ... d e ... d e ... e e a 1 ... f he ... a ... ec a ha ... e a i ... a ... ch
 a ... i b e he 1 1 1 a d i c 1 1 a ... c ... e ... f he ... a 1 g e a ... e . T h e ... e
 e ... ac ed , e ... e a 1 ... e e ... , acc ... d i g ... he e ... e 1 ... e a e i d e c e , ... e a i
 he d i c 1 1 a ... c ... e ... f he e a 1 g e a ... e . O c e he e ... e ... e e a 1
 1 e ... ac ed ... e a ... c e e d ... a ... i c a d a a a ... i ... c e a 1 ... he e a g g e 1 e
 ... e h d ... f fea ... e ... e e c 1 ... c ... d ... b e ... e d ... he e ... e ... e e a 1 ...

References

1. Morris, J., Coombes, K., Koomen, J., Baggerly, K., Kobayashi, R.: Feature extraction and quantification for mass spectrometry in biomedical applications using the mean spectrum. *Bioinformatics* (2005) Advanced publication.
2. Hastie, T., Tibshirani, R., Friedman, J.: *The Elements of Statistical Learning*. Springer (2001)
3. Prados, J., Kalousis, A., Sanchez, J.C., Allard, L., Carrette, O., Hilario, M.: Mining mass spectra for diagnosis and biomarker discovery of cerebral accidents. *Proteomics* **4** (2004) 2320–2332
4. Mallat, S.: *A wavelet tour of signal processing*. Academic Press (1999)
5. Petricoin, E., et al: Use of proteomic patterns in serum to identify ovarian cancer. *The Lancet* **395** (2002) 572–577
6. Petricoin, E., et al: Serum proteomic patterns for detection of prostate cancer. *Journal of the NCI* **94** (2002)
7. Witten, I., Frank, E.: *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann (1999)
8. Qu, Y., et al: Data reduction using a discrete wavelet transform in discriminant analysis of very high dimensional data. *Biometrics* **59** (2003) 143–151
9. Lee, K.R., Lin, X., Park, D., Eslava, S.: Megavariate data analysis of mass spectrometric proteomics data using latent variable projection method. *Proteomics* **3** (2003)
10. Zbigniew R. Struzik, A.S.: The haar wavelet transform in the time series similarity paradigm. In: *Principles of Data Mining and Knowledge Discovery, Third European Conference*, Springer (1999) 12–22

Numbers in Multi-relational Data Mining

Arno J. Knobbe^{1,2} and Eric K.Y. Ho¹

¹ Kiminkii, Postbus 171, NL-3990 DD, Houten, The Netherlands
{a.knobbe, e.ho}@kiminkii.com

² Utrecht University, P.O. box 80 089, NL-3508 TB Utrecht, The Netherlands

Abstract. Numeric data has traditionally received little attention in the field of Multi-Relational Data Mining (MRDM). It is often assumed that numeric data can simply be turned into symbolic data by means of discretisation. However, very few guidelines for successfully applying discretisation in MRDM exist. Furthermore, it is unclear whether the loss of information involved is negligible. In this paper, we consider different alternatives for dealing with numeric data in MRDM. Specifically, we analyse the adequacy of discretisation by performing a number of experiments with different existing discretisation approaches, and comparing the results with a procedure that handles numeric data dynamically. The discretisation procedures considered include an algorithm that is insensitive to the multi-relational structure of the data, and two algorithms that do involve this structure. With the empirical results thus obtained, we shed some light on the applicability of both dynamic and static procedures (discretisation), and give recommendations for when and how they can best be applied.

1 Introduction

Whereas numeric data is at the core of the majority of propositional Data Mining systems, it has been largely overlooked in Multi-Relational Data Mining (MRDM). Most MRDM systems assume that the data is a mixture of symbolic and structural data, and if the source database contains numbers, they will either have to be filtered out or pre-processed into symbolic values. Apart from historical reasons – symbolic representations are popular in the logical roots of MRDM –, the full treatment of numeric data comparable to propositional approaches is mostly ignored for reasons of simplicity and efficiency. MRDM is characterised by large hypothesis spaces, and the inclusion of continuous domains that offer a large range of (very similar) refinements is thought to make MRDM intractable. Most multi-relational systems rely on so-called discretisation procedures to reduce the continuous domains to more manageable symbolic domains of low cardinality, such that the search remains realistic. The resulting loss of precision is assumed to be negligible.

In this paper, we survey a number of existing approaches to dealing with numeric data in MRDM, with the aim of empirically determining the value of each of these approaches. These approaches include a number of pre-processing procedures suggested recently [6, 2], as well as one of the few MRDM algorithms that deal with numbers dynamically, developed by the authors of this paper [2, 4]. The discretisation procedures include a simple algorithm that considers each table in isolation, and discretises each numeric attribute on the basis of the distribution of its values,

regardless of any other tables connected to the current table. Two further discretisation procedures do involve the multi-relational structure of the database, and aim at finding good intervals, keeping in mind that the resulting symbolic attributes will be used in the context of the other tables in the database. The algorithm that deals with numbers dynamically does not require any pre-processing of the data. Rather than fixing a number of intervals prior to the analysis, it will consider the numeric data for a hypothesis at hand, and determine thresholds that are optimal for the given context. Especially at deeper levels of the search, where reasonably specific subgroups are considered, relevant thresholds will differ significantly from those determined on the whole dataset.

We test the four approaches experimentally on four well-known multi-relational datasets where numeric attributes play an important role: Mutagenesis (two varieties), Financial and Musk. With these experiments, we aim to shed some light on when and how each approach can best be applied. Furthermore, we hope to get some guidelines for important parameters of the discretisation procedures, such as the coarseness of the discretisation and the choice of representation. The experimental results are compared to those obtained on databases where all numeric information is removed, in order to get a baseline for the procedures that do (to some extent) involve the continuous domains.

2 Foundations

In the class of discrete patterns that we aim at (decision trees, rules, etc.), dealing with numeric data comes down to choosing numeric thresholds that form useful subgroups. Clearly, the distribution of numeric values, and how the target concept depends on this distribution is essential. In propositional data mining, choosing thresholds is fairly straightforward, as there is a one-to-one correspondence between occurring values and individuals. In MRDM however, we are dealing with non-determinate (i.e. one-to-many) relations between tables. In many cases, numeric attributes do not appear in the target table, and multiple values of the attribute are associated with a single structured individual. Whereas in propositional data mining, we can think of the whole database as a ‘cloud’ of points, in MRDM each individual forms a cloud. The majority of pattern languages in MRDM characterise such individuals by testing for the presence of values that exceed a given threshold. As the following lemma shows, only the largest and smallest values within each individual are relevant to include or exclude an individual on the basis of a single numeric test. Only these values will therefore be candidates for numeric thresholds.

Lemma 1. Let B be a bag of real numbers, and t some real, then

$$\begin{aligned} \exists v \in B: v \geq t & \text{ iff } \max(B) \geq t, \\ \exists v \in B: v \leq t & \text{ iff } \min(B) \leq t. \end{aligned}$$

Lemma 1 furthermore demonstrates that there is a difference between the set of thresholds appropriate for the \leq and the \geq operator. This means that any procedure that selects thresholds will have to be performed separately for each operator.

Choosing thresholds can roughly be done in two ways: dynamically and statically. A *dynamic* approach (see Section 3) considers the hypothesis at hand, and determines a collection of thresholds on the basis of the information contained in the individuals covered by the hypothesis in question. A *static* approach (see Section 4) on the other hand considers the entire database prior to analysis and determines a collection of thresholds once and for all. Typically these thresholds are then used to pre-process the data, replacing the numeric data with symbolic approximations. We refer to such a pre-processing step as *discretisation*. Clearly, a dynamic approach is preferable from an accuracy standpoint, as optimal thresholds are computed for the situation at hand. On the other hand, dynamic computation of thresholds makes algorithms more complex, and less efficient.

In the context of discretisation, we refer to numeric thresholds as *cut points*. A collection of $n-1$ cut points splits the continuous domain into n intervals. A group of values falling in a specific interval is referred to as a *bin*.

In MRDM, it makes sense to not just consider the available numeric values in the computation of cut-points, by also the multi-relational structure of the database. In general, a table is connected to other tables by associations, some of which may be non-determinate (a single record in one table corresponds to multiple records in another table). The effect of such associations is thus that records in a table can be divided into *groups*, depending on the relation to records in the associated table. Considering the multi-relational structure in the computation of cut points is hence tantamount to considering the numeric value, as well as the group the value belongs to. In the remainder of this paper, we refer to groups as the sets implied by this multi-relational structure.

3 Dynamic Handling of Numbers

An MRDM algorithm that handles numbers dynamically considers a range of cut points for a given numeric attribute, and determines how each of these tentative cut points influences the quality of a multi-relational hypothesis under consideration. As the optimal cut point depends on the current hypothesis, and many hypotheses are considered by an MRDM algorithm, the set of relevant cut points cannot be determined from the outset. Rather, we will have to consider the subgroup at hand, and query the database for a list of relevant cut points, and associated statistics.

In general, all values for the numeric attribute that occur in the individuals covered by the hypothesis at hand can act as candidate cut points. In theory, this set of values can be quite large, which can make the dynamic generation of cut points very inefficient. The MRDM system Safarii [2, 4] uses an approach that considers only a subset of these values, thus reducing some of the work. It relies on the observation from Lemma 1 that only the extreme values within a bag of numbers are relevant in order to test the presence of values above or below a certain cut point. Safarii uses a database primitive (a predefined query template) called NumericCrossTable [2] that selects the minimum (maximum) value within each individual covered by the current hypothesis, and then groups over these extreme values to produce the desired counts. We thus get a more reasonable number of candidate refinements.

Unfortunately it is still not realistic to continue the search on the basis of each of these refinements. Safarii therefore selects from the reduced set of candidate refinements only the optimal one for further examination. Because the operators \leq and \geq produce two different sets of candidate refinements, we essentially get two refinements per hypothesis and numeric attribute encountered. Note that keeping only the optimal refinements introduces a certain level of greediness into the algorithm.

4 Discretisation

In this section, we briefly outline the three methods for discretising numeric data to be used in our experiments. We refer to [3] for a full description. Conceptually, discretisation entails defining a number of consecutive intervals on the domain of a numeric attribute, and replacing this attribute with a nominal attribute that represents the interval values fall into. The three methods are identical in how numeric attributes are transformed based on the intervals defined. The essential difference between the methods lies in how the cut points between intervals are computed.

The first method presented computes a (user-determined) number of cut points based on the distribution of values of the numeric attribute. It ignores the fact that data in a particular table will generally be considered in the context of that in other tables. The remaining two methods do consider the multi-relational structure of the data, and compute cut points assuming that discretised values will be considered after joining with tables that are directly attached to the table at hand.

Because the numeric data typically appears in tables other than the target table, it is not always straightforward to assign a class (which is related to the target table) to the value. All three methods are therefore class-blind (or *unsupervised*): the methods do not consider a predefined target concept. As a result, the transformed data can be used on a range of class-definitions.

Equal Height Histogram. The first algorithm computes cut points regardless of any multi-relational structure. It simply considers every numeric attribute in every table in turn and replaces it by a nominal attribute that preserves as much of the information in the original attribute as possible. A collection of cut points is computed that produces bins of (approximately) equal size. Such a procedure is known as *equal interval frequency*, or *equal height histogram*, which is the term we will adopt.

Equal Weight Histogram. The second discretisation procedure involves an idea proposed by Van Laer et al. [6]. The algorithm considers not only the distribution of numeric values present, but also the groups they appear in. It is observed that larger groups have a larger impact on the choice of cut points because they have more contributing numeric values. In order to compensate for this, numeric values are weighted with the inverse of the size of the group they belong to. Rather than producing bins of equal size, we now compute cut points to obtain bins of equal weight.

Aggregated Equal Height Histogram. Like the EqualWeight algorithm, the AggregatedEqualHeight algorithm proposed in [2] takes the multi-relational structure of the database into account in the computation of the cut points. The algorithm is centred around the idea that not all values within a group are relevant when inquiring about the presence of numeric values above or below some threshold. As was outlined

in Section 2, it suffices to consider the minimum and maximum value within a group. The idea of the `AggregatedEqualHeight` algorithm is hence to take the minimum value per group and compute an equal height histogram on these values, in order to discretise all values. The process is then repeated for the maximum per group. We thus get two new attributes per numeric attribute.

Representation. In our discussion of the different discretisation procedures, we have assumed that the outcome is a collection of nominal attributes, where each value represents one of the computed intervals. In fact when we produce n nominal values, we do not only lose some amount of precision (which we assume to be minimal), but also the inherent order between intervals. Although the inability to handle ordered domains (numeric or ordinal) is part of our motivation for applying discretisation, we can choose a representation that preserves the order information without having to accommodate for it explicitly. This representation involves $n-1$ binary attributes per original numeric attribute, one for each cut point. Rather than representing each individual interval, the binary attributes represent overlapping intervals of increasing size. By adding such attributes as conjuncts to the hypothesis through repeated refinements, a range of intervals can be considered. A further advantage of this representation is that the accuracy is less sensitive to the number of intervals as the size of the intervals does not decrease with the number of intervals. An important disadvantage of this representation is the space it requires. Especially with larger numbers of intervals, having $n-1$ new binary attributes per original attribute can become prohibitive.

In our experiments, we will consider both the nominal and the binary representation, and compare the results to determine the optimal choice. We will refer to the latter representation as *cumulative binary*.

5 Experiments

Although we have multiple approaches to dealing with numeric data to test, we have chosen to apply a single mining algorithm. This allows us to sensibly compare results. The algorithm of choice is the Rule Discovery algorithm contained in the `Safarii MRDM` package produced by the authors [2, 4]. This algorithm produces a set of independent multi-relational rules. The algorithm includes the dynamic strategy for dealing with numbers described in Section 3. In order to test the discretisation procedures, we have pre-processed the different databases by generating the desired discretised attributes, and removing the original numeric attributes. The different discretisation procedures were implemented in the pre-processing companion to `Safarii`, known as `ProSafarii`.

Although a range of evaluation measures and search strategies is available in `Safarii`, we have opted for rules of high *novelty*, discovered by means of *beam search* (beam width 100, maximum depth 6). A time limit of 30 minutes per experiment was selected. The algorithm offers filtering of rules by means of a computed convex hull in ROC space [2]. The area under the ROC curve gives a good measure of the quality of the discovered rule set, as it is insensitive to copies or redundant combinations of rules. We will use this measure (values between 0.5 and 1) to compare results.

We will test the different algorithms on the following three well-known multi-relational databases:

- **Mutagenesis [5].** A database containing structural descriptions of molecules. We use two varieties, called B2 and B3. B2 contains symbolic and structural information as well as a single numeric attribute describing the charge of each atom. B3 contains two additional attributes on the molecule-level.
- **Financial [7, 2].** A database containing seven tables, describing various activities of customers of a Czech bank.
- **Musk [1].** A database describing 166 continuous features of different conformations molecules may appear in.

In [3] we present a detailed overview of the results obtained. We summarize the main conclusions in the paragraphs below.

Discretisation Procedures. Let us begin by considering how well the discretisation procedures perform. The table below summarises how often each procedure is involved in a win or a tie (no other procedure is superior). Procedures are compared per setting of the number of bins, in order to get comparable results. It turns out that AggregatedEqualHeight is clearly the best choice for Financial and Musk. Surprisingly, the propositional procedure EqualHeight performs quite well on Mutagenesis B2. The results for EqualHeight and EqualWeight on Mutagenesis B3 are virtually identical, which should come as no surprise, as this database contains two powerful attributes in the target table. The multi-relational data is mostly ignored.

In every case, the use of discretised attributes is better than not using the numeric information altogether, although in a few cases the advantage was minimal.

	EqualHeight	EqualWeight	AggregatedEqualHeight
Mutagenesis B2	62.5%	50.0%	37.5%
Mutagenesis B3	75.0%	87.5%	75.0%
Financial	0%	12.5%	87.5%
Musk	0%	25%	75.0%

Discretisation vs. Dynamic Handling. So can the discretisation procedures compete with the dynamic approach to numeric data, or is it always best to use the latter? In the table below, we compare the performance of the collection of discretisation procedures to dynamic handling of numbers. Each row shows in how many of the $3 \times 4 \times 2 = 24$ runs discretisation outperforms the dynamic approach. In the majority of cases, the dynamic approach outperforms the discretisation procedures, as was expected. However, for every database, there are a number of choices of algorithm, representation and number of bins, for which discretisation can compete, or even give slightly better results (see [3] for details).

If the set of cut points considered by the dynamic approach in theory is a superset of that considered by any discretisation procedure, how can we explain the moderate performance of the dynamic algorithm in such cases? The main reason is that the dynamic algorithm is more greedy than the discretisation procedures, because of the way numeric attributes are treated. Of the many refinements made possible by the numeric attribute, only the optimal pattern is kept for future refinements. Therefore, good rules involving two or more numeric conditions may be overlooked. On the

other hand, the nominal attributes resulting from discretisation produce a candidate for each occurring value, rather than only the optimal one. Because beam search allows several candidates to be considered, it may occur that sub-optimal initial choices may lead to optimal results in more complex rules.

	discretisation	dynamic
Mutagenesis B2	5	19
Mutagenesis B3	9	15
Financial	0	24
Musk	1	23

Choice of representation. The comparison between the two proposed representations is clear-cut: the cumulative binary representation generally gives the best results (see table below). The few cases where the nominal representation was (slightly) superior can be largely attributed to lower efficiency caused by the larger hypothesis space of the cumulative binary approach.

Although the cumulative binary representation is very desirable from an accuracy point of view, in terms of computing resources and disk space, the cumulative binary approach can become quite impractical, especially with many bins. Particularly in the Musk database, which contains 166 numeric attributes, several limits of the database technology used were encountered.

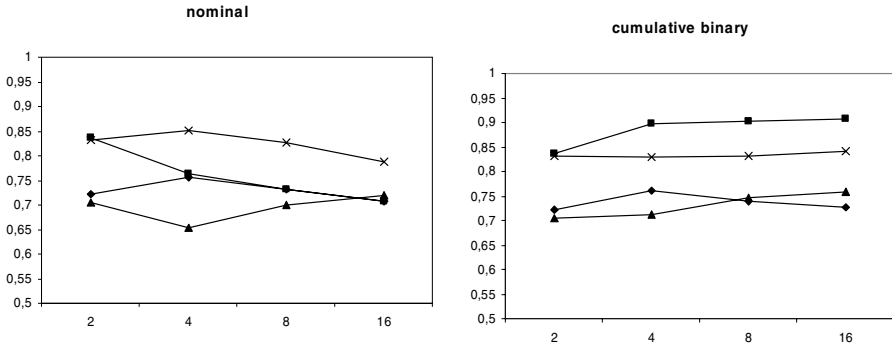
	nominal	cumulative binary	ties
Mutagenesis B2	3	5	4
Mutagenesis B3	0	9	3
Financial	2	6	4
Musk	5	4	7

Effect of Number of Bins. As has become clear, the number of bins is an important parameter of the discretisation procedures considered. Can we say something about the optimal value for this parameter? It turns out that the answer to this question depends on the choice of representation. Let us consider the cumulative binary representation. The performance roughly increases as more cut points are added (see the diagrams on the next page). This is because extra cut points just add extra opportunities for refinement and thus extra precision. The only exception to this rule is when severe time constraints are present. Because of the larger search space, there may be no time to reach the optimal result. For the nominal representation, there appears to be an optimal number of cut points that depends on specifics of the database in question. Having fewer cut points has a negative effect on the precision, whereas too many cut points results in rules of low support, because each nominal value only represents a small interval. For the Mutagenesis and Musk database, the optimal value is relatively low: between 2 and 4. The optimal value for Financial is less clear.

6 Conclusion

In general, we can say that the dynamic approach to dealing with numbers outperforms discretisation. This should come as no surprise, as the dynamic approach

is more precise in choosing the optimal numeric cut points. It is surprising however to observe that in some cases, it is possible to choose parameters and set up the discretisation process such that it is superior. Unfortunately, it is not immediately clear when faced with a new database what choice of algorithm, representation and



coarseness produces the desired result. Essentially, it is a matter of some experimentation to come up with the right settings. Even then, there is no guarantee that the extra effort of pre-processing the data provides a substantial improvement over the dynamic approach. Of course, when working with a purely symbolic MRDM system, discretisation is mandatory.

For discretisation, we recommend that the AggregatedEqualHeight procedure be tried first, as it has proven to give good results. It is worth the effort to consider EqualHeight as an alternative. The added value of the EqualWeight procedure over EqualHeight is negligible, and can therefore be ignored.

Our experimentation shows that in general, the simple nominal representation commonly used in MRDM projects is sub-optimal. Moreover, this representation is rather sensitive to the selected number of bins. In most cases the cumulative binary representation is preferable. This representation should be applied with as many bins as is realistic, given space and time limitations. Only when time restrictions can be expected to have a detrimental effect on the search depth, should lower numbers be considered.

References

1. Dietterich, T., Lathrop, R., Lozano-Pérez, T. *Solving the multiple-instance problem with axis-parallel rectangles*, Artificial Intelligence, 89(1-2):31-71, 1997
2. Knobbe, A.J. *Multi-Relational Data Mining*, Ph.D. dissertation, 2004, <http://www.kiminkii.com/thesis.pdf>
3. Knobbe, A.J., Ho, E.K.Y. *Numbers in Multi-Relational Data Mining*, 2005, <http://www.kiminkii.com/publications/pkdd2005long.pdf>
4. *Safarii, the Multi-Relational Data Mining engine*, Kiminkii, 2005, <http://www.kiminkii.com/safarii.html>
5. Srinivasan, A., Muggleton, S.H., Sternberg, M.J.E., King, R.D., *Theories for mutagenicity: A study in first-order and feature-based induction*, Artificial Intelligence, 85(1,2), 1996
6. Van Laer, W., De Raedt, L., Džeroski, S., *On multi-class problems and discretization in inductive logic programming*, In Proceedings ISMIS '97, LNAI 1325, Springer-Verlag, 1997
7. Workshop notes on Discovery Challenge PKDD '99, 1999

Testing Theories in Particle Physics Using Maximum Likelihood and Adaptive Bin Allocation

Bruce Knuteson¹ and Ricardo Vilalta²

¹ Laboratory for Nuclear Science, Massachusetts Institute of Technology,
77 Massachusetts Ave. Cambridge, MA 02139-4307, USA
knuteson@mit.edu

² Department of Computer Science, University of Houston,
4800 Calhoun Rd., Houston TX 77204-3010, USA
vilalta@cs.uh.edu

Abstract. We describe a methodology to assist scientists in quantifying the degree of evidence in favor of a new proposed theory compared to a standard baseline theory. The figure of merit is the log-likelihood ratio of the data given each theory. The novelty of the proposed mechanism lies in the likelihood estimations; the central idea is to adaptively allocate histogram bins that emphasize regions in the variable space where there is a clear difference in the predictions made by the two theories. We describe a software system that computes this figure of merit in the context of particle physics, and describe two examples conducted at the Tevatron Ring at the Fermi National Accelerator Laboratory. Results show how two proposed theories compare to the Standard Model and how the likelihood ratio varies as a function of a physical parameter (e.g., by varying the particle mass).

1 Introduction

Common to many scientific fields is the problem of comparing two or more competing theories based on a set of actual observations. In particle physics, for example, the behavior of Nature at small distance scales is currently well described by the Standard Model. But compelling arguments suggest the presence of new phenomena at distance scales now being experimentally probed, and there exists a long array of proposed extensions to the Standard Model.

The problem of assessing theories against observations can be solved in various ways. Some previous work bearing an artificial intelligence flavor has attempted to use observations to explain processes in both particle physics and astrophysics [4]. From a statistical view, a common solution is to use a maximum-likelihood approach [1,2], that selects the theory T maximizing $P(\mathcal{D}|T)$ (i.e., the conditional probability of a set of actual observations \mathcal{D} assuming T is true). Implicit to this methodology is the—often false—assumption that the form of the distributions characterizing the set of competing theories is known. In practice, a scientist suggests a new theory in the form of new equations or new parameters (e.g., new suggested mass for an elementary particle). In particle physics, a software is then used to simulate the response of the particle detector if the new proposed theory T were true, resulting in a data file made of Monte Carlo events from which one can estimate the true distribution characterizing T . At that point

one can compare how close T matches the actual observations (stored in \mathcal{D}) obtained from real particle colliders.

To estimate the true distribution of a theory T , we take the Monte Carlo data and follow a histogram approach [5]. We create a series of bins $\{b_k\}$ over the variable space and attempt to predict the number of events expected in every bin b_k if theory T were true. The novelty of our approach lies in the adaptive mechanism behind this bin allocation. Bins are selected to emphasize regions where the number of events predicted by T is significantly different from those predictions generated by competing theories, in a sense discovering regions in the variable space where a discrepancy among theories is evident.

This paper is organized as follows. Section 2 provides background information and notation. Section 3 provides a general description of the mechanism to compute likelihood ratios. Section 4 describes a solution to the problem of adaptive bin allocation. Section 5 reports on the experimental analysis. Lastly, Section 6 gives a summary and discusses future work.

2 Background Information and Notation

In modern particle accelerators, collisions of particles travelling at nearly the speed of light produce debris that is captured by signals from roughly one million channels of readout electronics. We call each collision an *event*. Substantial processing of the recorded signals leads to an identification of the different objects (e.g., electrons (e^\pm), muons (μ^\pm), taus (τ^\pm), photons (γ), jets (j), b -jets (b), neutrinos (ν), etc.) that have produced any particular cluster of energy in the detector. Each object is characterized by roughly three variables, corresponding to the three components of the particle's momentum. An event is represented as the composition of many objects, one for each object detected out of the collision. These kinematic variables can be usefully thought of as forming a *variable space*.

We store events recorded from real particle accelerators in a dataset $\mathcal{D} = \{\mathbf{e}_i\}$, where each event $\mathbf{e} = (a_1, a_2, \dots, a_n) \in A_1 \times A_2 \times \dots \times A_n$ is a variable vector characterizing the objects identified on a particular collision. We assume numeric variables only (i.e., $a_i \in \mathfrak{R}$) and that \mathcal{D} consists of independently and identically distributed (i.i.d.) events obtained according to a fixed but unknown joint probability distribution in the variable space.

We assume two additional datasets, D_n and D_s , made of discrete Monte Carlo events generated by a detector simulator designed to imitate the behavior of a real particle collider. The first dataset assumes the realization of a new proposed theory T_N ; the second dataset is generated under the assumption that the Standard Model T_S is true. Events follow the same representation on all three datasets.

3 Overview of Main Algorithm

In this section we provide a general description of our technique. To begin, assume a physicist puts forth an extension to the Standard Model through a new theory T_N . We define our metric of interest as follows:

$$\mathcal{L}(T_N) = \cdot_{10} \frac{P(\mathcal{D}|T_N)}{P(\mathcal{D}|T_S)} \tag{1}$$

where \mathcal{D} is the set of actual observations obtained from real particle colliders. Metric \mathcal{L} can be conveniently thought of as units of evidence for or against theory T_N . The main challenge behind the computation of \mathcal{L} lies in estimating the likelihoods $P(\mathcal{D}|\cdot)$. We explain each step next.

3.1 Partitioning Events into Final States

Each event (i.e., each particle collision) may result in the production of different objects, and thus it is appropriate to represent events differently. As an example, one class of events may result in the production of an electron; other events may result in the production of a muon. The first step consists of partitioning the set of events into subsets, where each subset comprises events that produced the same type of objects. This partitioning is orthogonal; each event is placed in one and only one subset, also called *final state*. Let m be the number of final states; the partitioning is done on all three datasets: $\mathcal{D} = \{\mathcal{D}_i\}_{i=1}^m$, $D_n = \{D_{ni}\}_{i=1}^m$, and $D_s = \{D_{si}\}_{i=1}^m$. Each particular set of subsets $\{\mathcal{D}_i, D_{ni}, D_{si}\}$ is represented using the same set of variables. Estimations obtained from each set of subsets are later combined into a single figure (Section 3.3).

3.2 Computation of Binned Likelihoods

The second step consists of estimating the likelihoods $P(\mathcal{D}|\cdot)$ adaptively by discovering regions in the variable space where there is a clear difference in the number of Monte Carlo event predictions made by T_N and T_S . Since we treat each subset of events (i.e., each final state) independently (Section 3.1), in this section we assume all calculations refer to a single final state (i.e. a single set of subsets of events $\{\mathcal{D}_i, D_{ni}, D_{si}\}$).

We begin by putting aside for a moment the real-collision dataset \mathcal{D}_i . The discrete Monte Carlo events predicted by T_N and T_S in datasets D_{ni} and D_{si} are used to construct smooth probability density estimates $P_i(\mathbf{e}|T_N)$ and $P_i(\mathbf{e}|T_S)$. Each density estimate assumes a mixture of Gaussian models:

$$P_i(\mathbf{e}|T) = P_i^T(\mathbf{e}) = \sum_{l=1}^r \alpha_l \phi(\mathbf{e}; \mu_l, \Sigma_l) \tag{2}$$

where r is the number of Gaussian models used to characterize the theory T under consideration. The mixing proportions α_l are such that $\sum_l \alpha_l = 1$, and $\phi(\cdot)$ is a multivariate normal density function:

$$\phi(\mathbf{e}; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{e} - \mu)^t \Sigma^{-1}(\mathbf{x} - \mu)} \tag{3}$$

where \mathbf{e} and μ are d -component vectors, and $|\Sigma|$ and Σ^{-1} are the determinant and inverse of the covariance matrix.

At this point we could follow the traditional approach to Maximum Likelihood estimation by using the real-collision dataset \mathcal{D}_i and the above probability density estimates:

$$P(\mathcal{D}_i|T) = \prod_j P_i(\mathbf{e}_j|T) = \prod_j P_i^T(\mathbf{e}_j) \quad (4)$$

where T takes on T_N or T_S and the index j goes along the events in \mathcal{D}_i .

The densities $P_i(\mathbf{e}|T)$ can in principle be used to compute an unbinned likelihood ratio. But in practice, this ratio can suffer from systematic dependence on the details of the smoothing procedure. Over-smoothed densities cause a bias in favor of distributions with narrow Gaussians, while the use of under-smoothed densities cause undesired dependence on small data irregularities. The calculation of a binned likelihood ratio in the resulting discriminant reduces the dependence on the smoothing procedure, and has the additional advantage that it can be used directly to highlight regions in the variable space where predictions from the two competing theories T_N and T_S differ significantly. We thus propose to follow a histogram technique [5] as follows.

Constructing a Binned Histogram

We begin by defining the following discriminant function:

$$D(\mathbf{e}) = \frac{P_i(\mathbf{e}|T_N)}{P_i(\mathbf{e}|T_N) + P_i(\mathbf{e}|T_S)} \quad (5)$$

The discriminant function D takes on values between zero and unity, approaching zero in regions in which the number of events predicted by the Standard Model T_S greatly exceeds the number of events predicted by the new proposed theory T_N , and approaching unity in regions in which the number of events predicted by T_N greatly exceeds the number of events predicted by T_S . We employ function D for efficiency reasons: it captures how the predictions of T_N and T_S vary in a single dimension.

We use D to adaptively construct a binned histogram. We compute the value of the discriminant D at the position of each Monte Carlo event predicted by T_N (i.e., every event contained in D_n) and T_S (i.e. every event contained in D_s). The resulting distributions in D are then divided into a set of bins that maximize an optimization function. This is where our adaptive bin allocation strategy technique is invoked (explained in detail in Section 4). The result is a set of bins that best differentiate the predictions made by T_N and T_S . The output of the Adaptive-Bin-Allocation algorithm is an estimation of the conditional probability $P(\mathcal{D}_i|T)$.

As an illustration, Figure 1 (left) shows the resulting binned histogram in D for a real scenario with a final state e^+e^- (i.e., electron and positron). The Adaptive-Bin-Allocation algorithm chooses to consider only two bins, placing a bin edge at $D = 0.4$. Note events from T_S (line L2) tend to lie at values for which $D(\mathbf{e})$ is small, and events from T_N (line L3) tend to lie at values for which $D(\mathbf{e})$ is large.

Figure 1 (right) shows how the two bins in the discriminant map back onto the original variable space defined on $m_{e^+e^-}$ (the invariant mass of the electron positron pair), and positron pseudorapidity. The dark region corresponds to points \mathbf{e} in the variable space for which $D(\mathbf{e}) < 0.4$; similarly the light region corresponds to points \mathbf{e}

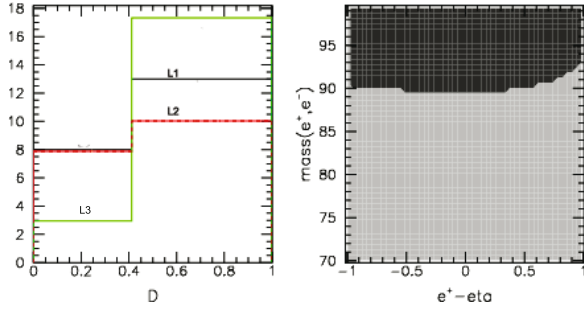


Fig. 1. (Left) The optimally-binned histogram of the discriminant D for the predictions of T_S (line L2), T_N (line L3), and real data \mathcal{D} (line L1). (Right) The mapping of the bins in D back into regions in the original variable space. The dark region corresponds to points \mathbf{e} in the variable space for which $D(\mathbf{e}) < \theta$; the light region corresponds to points \mathbf{e} in the variable space for which $D(\mathbf{e}) > \theta$ ($\theta = 0.4$).

for which $D(\mathbf{e}) > 0.4$. Each region is assigned a binned probability (Section 4); all probabilities are then combined into a final state probability $P(\mathcal{D}_i|T)$.

3.3 Combining Likelihoods and Incorporating Systematic Errors

Once we come up with an estimation of $P(\mathcal{D}_i|T)$, the next step consists of combining all probabilities from individual final states into a single probability for the entire experiment through the product $P(\mathcal{D}|T) = \prod_i P(\mathcal{D}_i|T)$, where T takes on T_N or T_S and the index i goes along all final states. As a side note, a single particle accelerator has normally several experiments running that can also be combined through such products.

Finally, systematic uncertainties are introduced into the analysis to reflect possible imperfections in the modelling of the response of the physical detector. There are usually roughly one dozen sources of systematic error, ranging from possible systematic bias in the measurements of particle energies to an uncertainty in the total amount of data collected.

4 Adaptive Bin Allocation

We now explain in detail our approach to estimate the likelihood $P(\mathcal{D}_i|T)$ for a particular final state. To begin, assume we have already computed the value of the discriminant D at the position of each Monte Carlo event predicted by T_N and T_S (Section 3.2), and decided on a particular form of binning that partitions D into a set of bins $\{b_k\}$. Let $\mu_{k|T}$ be the number of events expected in bin k if theory T is true¹. Often in the physical sciences the distribution of counts in each bin is Poisson; this is assumed in what follows. The probability of observing λ_k events in a particular bin k is defined as:

¹ Recall T is either the new theory T_N or the Standard Model T_S .

$$P(\lambda_k|T) = \frac{e^{-\mu_k|T} \mu_k|T^{\lambda_k}}{\lambda_k!} \tag{6}$$

Now, the probability of observing the real data D_i assuming the correctness of T and neglecting correlated uncertainties among the predictions of T in each bin, is simply:

$$P(D_i|T) = \prod_k P(\lambda_k|T) \tag{7}$$

where the index k runs along the bins and λ_k is the number of events observed in the real data D_i within bin k .

The question we now pose is how should the bins be chosen? Many finely spaced bins allow finer sampling of differences between T_N and T_S , but introduce a larger uncertainty in the prediction within each bin (i.e., the difference in the events predicted by T_N and T_S under finely spaced bin comes with low confidence levels). On the other hand, a few coarsely spaced bins allow only coarse sampling of the distributions predicted by T_N and T_S , but the predictions within each bin are more robust. The question at hand is not only how many bins to use, but also where to place their edges along the discriminant D [3].

4.1 Searching the Space of Binnings

In selecting an optimal binning we focus our analysis on the two theories T_N and T_S exclusively (choosing a set of optimal bins is independent of the real data used for theory validation). Our goal is to produce a set of bins $\{b_k\}$ that maximize the difference in predictions between the two theories. We start by defining an optimization function over the space of binnings. We merit partitions that enhance the expected evidence in favor of T_N , $\mathcal{E}(T_N)$, if T_N is correct, plus the expected evidence in favor of T_S , $\mathcal{E}(T_S)$, if T_S is correct. Given a particular set of bins, $\{b_k\}_{k=1}^v$, the proposed optimization function is defined as follows:

$$\mathcal{O}(\{b_k\}) = \mathcal{E}(T_N, \{b_k\}) + \mathcal{E}(T_S, \{b_k\}) \tag{8}$$

The evidence for each theory is as follows:

$$\mathcal{E}(T_N, \{b_k\}) = \sum_{\lambda_1} \sum_{\lambda_2} \cdots \sum_{\lambda_v} \left(\prod_k P(\lambda_k|T_N) \right) \times \cdot g_{10} \left(\frac{\prod_k P(\lambda_k|T_N)}{\prod_k P(\lambda_k|T_S)} \right) \tag{9}$$

and similarly for $\mathcal{E}(T_S, \{b_k\})$. Each summation on the left varies over the range $[0, \infty]$. The evidence for each theory has a straightforward interpretation. Recall that $\prod_k P(\lambda_k|T) = P(D_i|T)$ and therefore each evidence \mathcal{E} is the relative entropy of the data likelihoods (if $\cdot g_{10}$ is replaced with $\cdot g_2$), averaged over all possible outcomes on the number of real events observed on each bin. The two components in equation 8 are necessary because relative entropy is not symmetric. The representation for \mathcal{O} can be simplified as follows:

Algorithm 1: Adaptive-Bin-Allocation

Input: $D, \tilde{D}_{n_i}, \tilde{D}_{s_i}$

Output: Set of bins $\{b_k\}$

ALLOCATE-BINS($D, \tilde{D}_{n_i}, \tilde{D}_{s_i}$)

- (1) Evaluate D at each discrete Monte Carlo event in \tilde{D}_{n_i} and \tilde{D}_{s_i} .
- (2) Estimate probability densities $f(\mu_{k|T})$ for $T = T_N$ and $T = T_S$.
- (3) Initialize set of bins $\{b_0\}$, where b_0 covers the entire domain of D .
- (4) **repeat**
- (5) Search for a cut point c over D that maximizes function \mathcal{O} .
- (6) Replace the bin b_k where c falls with the two corresponding new bins.
- (7) **until** The value o^* maximizing $\mathcal{O}(\cdot)$ is such that $o^* < \epsilon$
- (8) **end**
- (9) **return** $\{b_k\}$

Fig. 2. Steps to generate a set of bins that maximize the distance between the events predicted by theory T_N and theory T_S

$$\mathcal{O}(\{b_k\}) = \sum_k \sum_{\lambda_k} (\mathbb{P}(\lambda_k|T_N) - \mathbb{P}(\lambda_k|T_S)) \times (\cdot_{g_{10}} \mathbb{P}(\lambda_k|T_N) - \cdot_{g_{10}} \mathbb{P}(\lambda_k|T_S)), \tag{10}$$

In practice one cannot evaluate \mathcal{O} by trying all possible combinations in the number of real events observed on each bin. Instead we estimate the expected number of events in bin k if theory T is true, $\mu_{k|T}$, and consider $\pm s$ standard deviations (s is user-defined) around that expectation, which can be quickly evaluated with arbitrary accuracy by explicitly computing the sum for those bins with expectation $\mu_{k|T} \leq 25$ and using a gaussian approximation for those bins with expectation $\mu_{k|T} > 25$.

Although in principle maximizing \mathcal{O} requires optimizing the positions of all bin edges simultaneously, in practice it is convenient to choose the bin edges sequentially. Starting with a single bin encompassing all points, this bin is split into two bins at a location chosen to maximize \mathcal{O} . At the next iteration, a new split is made that improves \mathcal{O} . The algorithm continues iteratively until further division results in negligible or negative change in \mathcal{O} . Figure 2 (Algo. 1) illustrates the mechanism behind the binning technique. The complexity of the algorithm is linear in the size of the input space (i.e., in the size of the two datasets D_{n_i} and D_{s_i}).

4.2 Example with Gaussians of Varying Width

To illustrate the mechanism behind the bin-allocation mechanism, assume a scenario with two Gaussian distributions of different widths over a variable x . Figure 3(left) shows the true (but unknown) distributions $f_1(x)$ and $f_2(x)$, where $f_i(x) = \frac{n}{\sqrt{2\pi}\sigma_i} e^{-(x-\mu)^2/2\sigma_i^2}$ with $i = \{1, 2\}$ and parameter values $n = 100, \mu = 25, \sigma_1 = 5,$ and $\sigma_2 = 8$. The units on the vertical axis are the number of events expected in the data per unit x . We used one thousand points randomly drawn from $f_1(x)$ and from $f_2(x)$. These points are shown in the histogram in Fig. 3(right), in bins of unit width in x . The algorithm proceeds to find edges sequentially before halting, achieving a final

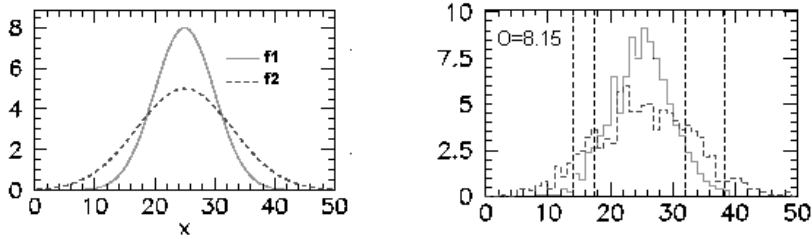


Fig. 3. (Left) Two Gaussian distributions, f_1 and f_2 , with same mean but different variance. (Right) The bin-allocation mechanism identifies those regions where f_1 and f_2 cross.

figure of merit. The resulting bins are concentrated in the regions $x \approx 20$ and $x \approx 30$, where $f_1(x)$ and $f_2(x)$ cross.

5 Experiments

We describe two examples conducted at the Tevatron ring at the Fermi National Accelerator Laboratory in Chicago, Illinois. The accelerator collides protons and anti-protons at center of mass energies of 1960 GeV (i.e., giga electron volts). A typical real-collision dataset of this collider is made of about 100 thousand events.

We divide each of the Monte Carlo data sets D_n , and D_s into three equal-size subsets. The first subset is used to compute the probability densities $P_i(\mathbf{e}|T_N)$, $P_i(\mathbf{e}|T_S)$ (Section 3.2); the second subset is used to run the adaptive bin-allocation mechanism (Section 4); the last subset is used to estimate the figure of merit $\mathcal{L}(T_N) = g_{10} \frac{P(\mathcal{D}|T_N)}{P(\mathcal{D}|T_S)}$ (Section 3). Each experiment produces several hundreds of final states. The running time for each experiment was approximately one hour on a Linux machine with a Pentium 3 processor and 1 GB of memory.

Searching for Leptoquark Pair Production. The first experiment is motivated by a search for leptoquark pair production as a function of assumed leptoquark mass. We show how a theory that advocates leptoquarks with small masses –that if true would result in an abundance of these particles compared to their heavier counterparts– is actually disfavored by real data. Figure 4 (left) shows the log likelihood ratio $\mathcal{L}(T_N)$ (equation 1) for different leptoquark masses. Units on the horizontal axis are GeV. The new proposed theory is disfavored by the data for small mass values, but becomes identical to the Standard Model for large mass values. Figure 4 (second left) shows the posterior distribution $p(M_{LQ}|\mathcal{D})$ obtained from a flat prior and the likelihood on the left.

Searching for a Heavy Z' Particle. The second experiment is similar in spirit to the previous one. Figure 4(third from left) shows a search for a heavy Z' as a function of assumed Z' mass. Z' s with small masses, which would be more copiously produced in the Tevatron than their heavier counterparts, are disfavored by the data. The posterior probability $p(m_{Z'}|\mathcal{D})$ flattens out beyond $m_{Z'} \approx 250$ GeV (Figure 4, right), indicating that the data is insufficiently sensitive to provide evidence for or against Z' s at this mass.

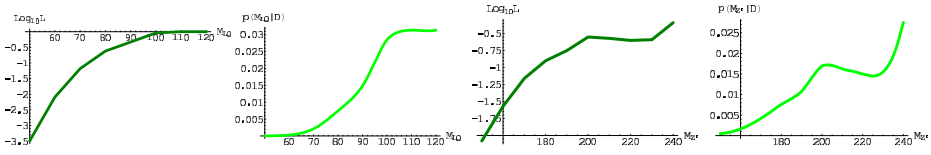


Fig. 4. (left) The log likelihood ratio $\mathcal{L}(T_N)$ (equation 1) for different leptokuark masses. (second left) The posterior distribution $p(M_{LQ}|\mathcal{D})$ obtained from a flat prior and the likelihood on the left. (third left) The log likelihood ratio for different Z' masses. (right) The posterior probability $p(m_{Z'}|\mathcal{D})$ flattens out beyond $m_{Z'} \approx 250$ GeV. Units on the horizontal axis are GeV.

6 Conclusions and Future Work

This paper describes an approach to quantify the degree of evidence in favor of a new proposed theory compared to a standard baseline theory. The mechanism adaptively allocates histogram bins that emphasize regions in the variable space where there is a clear difference in the predictions made by the two theories. The proposed mechanism carries two important benefits: 1) it simplifies substantially the current time needed to assess the value of new theories, and 2) it can be used to assess a family of theories by varying a particular parameter of interest (e.g., particle mass).

We expect the procedure outlined here to have widespread application. The calculation of likelihood ratios is common practice in the physical and social sciences; the main algorithm can be easily adapted to problems stemming from other scientific fields. One barrier lies in generating Monte Carlo data to model a theory distribution. Particle physicists have invested huge amounts of effort in producing a detector simulator designed to imitate the behavior of real particle colliders.

Acknowledgments. This material is based upon work supported by the National Science Foundation under Grants no. IIS-431130 and IIS-448542.

References

1. Duda R. O., Hart P. E., Stork D. G.: Pattern Classification. John Wiley Ed. 2nd Edition (2001).
2. Hastie T., Tibshirani R., Friedman J.: The Elements of Statistical Learning. Springer-Verlag Ed. (2001).
3. Knuteson, Bruce: Systematic Analysis of HEP collider data. Conference on Statistical Problems in Particle Physics, Astrophysics, and Cosmology. Stanford CA.
4. Kocabas S., Langley P.: An Integrated Framework for Extended Discovery in Particle Physics. Proceedings of the 4th International Conference on Discovery Science, pp. 182-195. Springer Verlag. (2001).
5. Scott D.W.: Multivariate Density Estimation: Theory, Practice, and Visualization. Wiley Series in Probability and Statistics Ed. (1992).

Improved Naive Bayes for Extremely Skewed Misclassification Costs

Alejandro K. C. Abdur Chahid

AOL, Inc., 44900 Prentice Drive, Dulles VA 20166, USA
{arkolcz, cabdur}@aol.com

Abstract. Naive Bayes has been an effective and important classifier in the text categorization domain despite violations of its underlying assumptions. Although quite accurate, it tends to provide poor estimates of the posterior class probabilities, which hampers its application in the cost-sensitive context. The apparent high confidence with which certain errors are made is particularly problematic when misclassification costs are highly skewed, since conservative setting of the decision threshold may greatly decrease the classifier utility. We propose an extension of the Naive Bayes algorithm aiming to discount the confidence with which errors are made. The approach is based on measuring the amount of change to feature distribution necessary to reverse the initial classifier decision and can be implemented efficiently without over-complicating the process of Naive Bayes induction. In experiments with three benchmark document collections, the decision-reversal Naive Bayes is demonstrated to substantially improve over the popular multinomial version of the Naive Bayes algorithm, in some cases performing more than 40% better.

1 Introduction

In the binary classification problem, the high decision cost of the high misclassification rate can be reduced if the classifier is able to estimate the feature probabilities more accurately. The high decision cost of the high misclassification rate can be reduced if the classifier is able to estimate the feature probabilities more accurately. The high decision cost of the high misclassification rate can be reduced if the classifier is able to estimate the feature probabilities more accurately.

In this paper, we extend the decision-reversal Naive Bayes classifier to handle skewed misclassification costs. Our extension is based on measuring the amount of change to feature distribution necessary to reverse the initial classifier decision and can be implemented efficiently without over-complicating the process of Naive Bayes induction. In experiments with three benchmark document collections, the decision-reversal Naive Bayes is demonstrated to substantially improve over the popular multinomial version of the Naive Bayes algorithm, in some cases performing more than 40% better.

3.2 Overconfidence in Decision Making

In this paper, we have shown that Naive Bayes, which is generally considered to be a good classifier, can be overconfident in its decision making. This is especially true for documents that are misclassified. We have shown that the confidence of the classifier is high for documents that are misclassified, even when the training data contains much of related content. This is a problem because it means that the classifier is not only wrong, but also very sure about its wrong decision. This is a problem because it means that the classifier is not only wrong, but also very sure about its wrong decision. This is a problem because it means that the classifier is not only wrong, but also very sure about its wrong decision.

Figure 1 shows the scores of Naive Bayes misclassifications (at default decision threshold) vs. maximum training-set document frequency for features belonging to the misclassified documents (for the collections of: Reuters-21578, 20-Newsgroups and TREC-AP). Misclassifications of documents falling into sparsely populated regions are likely to be made with higher confidence (signified by high absolute score values) than those made for documents for which the training data contained much of related content.

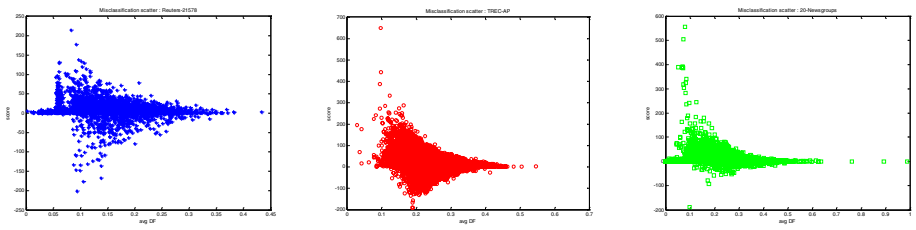


Fig. 1. Scores of Naive Bayes misclassifications (at default decision threshold) vs. maximum training-set document frequency for features belonging to the misclassified documents (for the collections of: Reuters-21578, 20-Newsgroups and TREC-AP). Misclassifications of documents falling into sparsely populated regions are likely to be made with higher confidence (signified by high absolute score values) than those made for documents for which the training data contained much of related content.

This can be explained by the fact that documents with high maximum training-set document frequency for features belonging to the misclassified documents are likely to be made with higher confidence (signified by high absolute score values) than those made for documents for which the training data contained much of related content.

In [6] it is argued that the high confidence of the classifier is a result of the fact that the classifier is not only wrong, but also very sure about its wrong decision. This is a problem because it means that the classifier is not only wrong, but also very sure about its wrong decision. This is a problem because it means that the classifier is not only wrong, but also very sure about its wrong decision.

additionally, we assume that $P(C|x) > 0$ for all $x \in \mathcal{X}$. We have $P(C|x) = \frac{P(C, x)}{P(x)}$ and $R(C|x) = \frac{P(C, x)}{P(C)}$. The joint probability $P(C, x)$ can be written as $P(C, x) = P(x|C) \cdot P(C)$. The joint probability $P(x|C)$ can be written as $P(x|C) = \frac{P(x, C)}{P(C)}$. The joint probability $P(x, C)$ can be written as $P(x, C) = P(x|C) \cdot P(C)$. The joint probability $P(x, C)$ can be written as $P(x, C) = P(x|C) \cdot P(C)$.

$$\hat{P}(C|x) = P(C|x) \cdot R(C|x) \tag{2}$$

where $R(C|x)$ is the reliability of the classifier C for the class C given the input x .

4 Changing Naive Bayes' Mind: A New Reliability Measure

The goal of this paper is to propose a new reliability measure for the Naive Bayes classifier. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features.

The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features.

The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features. The Naive Bayes classifier is a probabilistic classifier that uses Bayes' theorem to calculate the probability of each class given the input features.

$$\log P(x|C) - \log P(x|\bar{C}) = \log \tilde{P}(x|\bar{C}) - \log \tilde{P}(x|C) - score \tag{3}$$

where $score$ is the log-likelihood ratio of the two classes, i.e., $\log \tilde{P}(x|C) - \log \tilde{P}(x|\bar{C})$.

As the available data has been subject to the expected change in the distribution. Hence we consider the Kullback-Leibler (KL) divergence, i.e.,

$$rdist(x) = KL \left(P(x|\bar{C}), \tilde{P}(x|\bar{C}) \right) = \sum_{x_i} P(x_i|\bar{C}) \cdot \log \frac{P(x_i|\bar{C})}{\tilde{P}(x_i|\bar{C})} \quad (4)$$

Once the KL divergence (4) is computed, a straightforward combination of the distribution can be achieved. (2) the original distribution of the predicted class. This has been addressed in the context of the KL divergence, in [6] and [7]. Hence the distribution of the predicted class is given by $R(C|x) : rdist(x) \rightarrow [0, 1]$, by adding a bias to the original distribution of Naive Bayes. The original distribution of Naive Bayes is denoted by $P(C|x)$. The modified distribution of (4) is denoted by $R(C|x)$ and $P(C|x)^1$. Given that the classifier is a naive Bayes classifier, the original distribution of Naive Bayes is given by $P(C|x) = \frac{score(x)}{score(x) + 1}$, where $score(x)$ is the score of the classifier. If $score(x) = 1$, then $P(C|x) = 1$. If $score(x) = 0$, then $P(C|x) = 0$. Hence the bias is added to the original distribution of Naive Bayes. The effect of the bias is to reduce the distribution of Naive Bayes. The effect of the bias is to reduce the distribution of Naive Bayes. The effect of the bias is to reduce the distribution of Naive Bayes.

$$\widehat{score}(x) = score(x) \cdot rdist(x) \quad (5)$$

Once the bias is added to the original distribution of Naive Bayes, the effect of the bias is to reduce the distribution of Naive Bayes. The effect of the bias is to reduce the distribution of Naive Bayes. The effect of the bias is to reduce the distribution of Naive Bayes.

$$\widehat{score}(x) = score(x) \cdot e^{-\gamma \cdot rdist(x)} \quad (6)$$

where γ is a parameter.

5 Experimental Setup

The experiments described below are carried out on the dataset where the classifier achieves 100% accuracy for the age class. A high degree of accuracy is achieved by the classifier, i.e., the classifier is able to correctly classify the age class. The accuracy of the classifier is high, i.e., the classifier is able to correctly classify the age class. The accuracy of the classifier is high, i.e., the classifier is able to correctly classify the age class.

We compare the proposed method with the Naive Bayes (a-b) and NB-KL in the following:

- NB: Using the original Naive Bayes (baseline).
- NB-Trans: Kullback-Leibler divergence based Naive Bayes [7] (highly effective in the case of extremely skewed distributions).

¹ In fact [7] does it directly by substituting the posterior estimate of $P(C|x)$ with $prec \cdot R(C|x)$, where $prec$ refers to the overall precision of the classifier.

Table 1. Steps involved in the decision-reversal Naive Bayes. The most computation-ally expensive part is step 2, in which one needs to estimate how many corrective events need to take place before the initial decision of the classifier is changed. A naive implementation would keep on generating such events and updating the model, but since in some cases the number of events may be on the order of hundreds or more, this would add significantly to the evaluation time. Instead, we treat the score as a function of the corrective event count α and identify the zero-crossing of score(α). In our implementation of the Newton method, usually only 1–7 iterations are needed.

<p>Algorithm</p> <ol style="list-style-type: none"> 1. Classify input x using a trained NB model. 2. Estimate the multiplicity α with which x needs to be added to the opposite class to achieve decision reversal. 3. Measure the KL divergence (eq.(4)) between the original and the perturbed distribution of features for the class opposite to the one originally predicted. 4. Modulate the original score (eq.(5) or (6)).
--

Multi-class models are evaluated as a function of class accuracy, which is calculated as the average of the per-class accuracies over all classes, i.e., $\frac{1}{C} \sum_{c \in \mathcal{C}} \text{acc}_c$. The per-class accuracy is defined as the fraction of correctly classified instances for each class c .

5.1 Data Sets

We chose three datasets which have been extensively used in the text classification literature. In each case the dataset is split (1/3 held out as a development set, the rest for training) into a training and a development set as follows:

- Reuters-21578 (101 categories, 10,724 documents): We used the standard `mod_lapte` split of the data.
- 20 Newsgroups (20 categories, 19,997 documents): A random sample of 2/3 of the dataset was chosen for training with the remaining documents used for testing.
- TREC-AP (20 categories, 209,783 documents): The training/test split described in [?] was used.

Features were extracted by using a bag-of-words model, where each word is represented by its frequency in the document. We used the standard `tf-idf` normalization. In addition, we used a bag-of-words model to extract features from the document titles.

Table 2. Macro-averaged classification performance (non-target specificity) captured at the point of perfect target recall. The decision-reversal variant of Naive Bayes consistently outperformed the baseline, while the transductive method consistently underperformed in all three cases.

Dataset	NB	NB-Trans	NB-KL	$\frac{\Delta(\text{NB-KL}-\text{NB})}{\text{NB}}$ [%]
Reuters-21578	0.4743	0.3070	0.6693	41
20 Newsgroups	0.4033	0.3297	0.5379	33
TREC-AP	0.5004	0.1871	0.5954	19

In a γ -case (e.g., the feature set reduced to the top 5,000 attributes), the highest accuracy of the MI baseline (MI) between the feature labels and the class labels is achieved. The algorithm

6 Results

Table 2 shows the results. For all three datasets, NB-KL provided a substantial improvement over the baseline NB. The standard exponential discounting [7] generally outperformed the baseline NB. With high discounting, however, the performance of NB-KL drops significantly. To achieve high specificity at a 100% target recall, the exponential discounting of the age class is essential. In addition, the high specificity of the age class is essential for the high specificity. In each case, the decision reversal variant of NB consistently outperforms the baseline NB. The standard exponential discounting of the age class is essential for the high specificity. In each case, the decision reversal variant of NB consistently outperforms the baseline NB. In addition, the high specificity of the age class is essential for the high specificity.

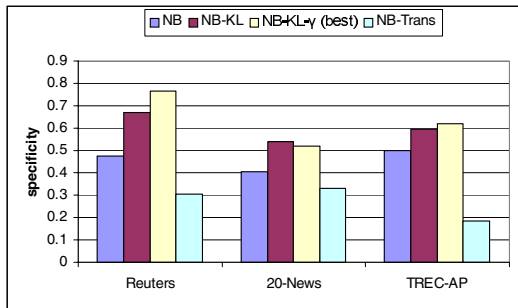


Fig. 2. at the point of perfect target recall. The results for best parameter settings in (6) are compared to the baseline NB, NB-Trans and the default settings of NB-KL. In the case of Reuters-21578 and TREC-AP exponential discounting results in substantial increase in specificity. For 20-Newsgroups, however, the original formulation of NB-KL works better.

Clustering and Prediction of Mobile User Routes from Cellular Data

Kari Laasonen

Basic Research Unit, Helsinki Institute for Information Technology,
Department of Computer Science, University of Helsinki
Kari.Laasonen@cs.Helsinki.FI

Abstract. Location-awareness and prediction of future locations is an important problem in pervasive and mobile computing. In cellular systems (e.g., GSM) the serving cell is easily available as an indication of the user location, without any additional hardware or network services. With this location data and other context variables we can determine places that are important to the user, such as work and home. We devise online algorithms that learn routes between important locations and predict the next location when the user is moving. We incrementally build clusters of cell sequences to represent physical routes. Predictions are based on destination probabilities derived from these clusters. Other context variables such as the current time can be integrated into the model. We evaluate the model with real location data, and show that it achieves good prediction accuracy with relatively little memory, making the algorithms suitable for online use in mobile environments.

1 Introduction

Location-awareness has a great potential in pervasive computing. Several applications have been developed that use location information for predicting the user's future location. The user's location is usually obtained from the serving cell of the mobile phone. The location data is often used to determine important locations, such as home and work, and to predict the user's future location. The prediction is based on the user's past location data. The prediction is often done by clustering the location data into clusters, and then predicting the next location based on the clusters. The prediction is often done by using a Markov chain model, where the states are the clusters of location data. The prediction is often done by using a hidden Markov model, where the states are the clusters of location data. The prediction is often done by using a neural network, where the input is the current location and the output is the predicted next location.

This paper describes a new method for predicting the user's future location. The method is based on clustering the location data into clusters, and then predicting the next location based on the clusters. The prediction is done by using a Markov chain model, where the states are the clusters of location data. The prediction is often done by using a hidden Markov model, where the states are the clusters of location data. The prediction is often done by using a neural network, where the input is the current location and the output is the predicted next location.

M...e... de... i... g... ca... a... e... GPS c...
di... da... [1,2,3]. H...e..., GPS ca... be... ba... a... de...
... g... had... g... GPS... e... a... h...e... b... i... a... b...
... h...e... A h... a... d... S...e... [1] c... e... c... di... e... da... a... i... f...e... ca... .., b...
... e... e... i... e... ca... be... ed... e... [5]. A... e... a... e... h... d... f...e... d... i... c... f...
f...e... ca... .. i... c... de... .. e... c... d... de... Ma... .. de... [1], a... d... Ba... e... ca...
... i... e... [2,6].

O... da... a... e... h...e... f...e... f...a... e... e... c... e... f... ce... i... de... i... e... . A... i... e... e... i... g... a...
... ach... c... e... i... g... e... e... ce... i... h... h... b...a... i... i... c... .. e... [7]. S... ch... e... h... d...
... f...e... f... a... e... e... i... e... .. ch... e... .. a... d... .. ce... i... g... ca... a... c... i... .. be... f...e...
... b... i... e... h...e... .

2 Problem Setting

A GSM h...e... c... .. i... ca... e... .. e... h...e... a... i... h... a... b... e... a... i... . I... a... g... e...
... ca... .. h...e... e... a... b... e... e... a... b... a... e... a... i... .. h...e... a... d... i... g... a... .. e... a... ch... e... h...e... h...e...
The h...e... ch... .. e... e... f... h...e... .. a... d... .. i... c... h...e... .. a... .. a... e... .. e... .. e... .. a... e... .. b... a...
... a... i... .. a... .. e... e... d... . A... .. i... .. h...e... a... e... a... c... e... e... d... b... a... i... g... e... b... a... e... a... i... .. ; h...e...
... a... h...e... h...e... i... .. i... .. e... .. ce... .. , e... .. e... a... .. h...e... h...e... i... .. h...e... a... e... a... f... h...e...
c... .. e... .. d... i... g... b... a... e... a... i... .. .

I... h...i... a... e... .. e... .. i... h... GSM ce... .. da... a... f... .. a... .. b... e... .. f... .. e... a... .. . M... b... i... e...
... h...e... a... e... b... i... .. i... .. a... d... .. ce... .. a... .. e... .. e... .. a... .. e... .. e... .. h...e... e... . S... i... ce... ..
... .. e... .. a... .. e... .. e... .. a... .. e... .. i... c... e... i... f... a... .. c... .. e... .. a... .. e... .. i... .. e... .. d... .. da... .. a... .. g... a... .. h...e... i... g... i... e... a...
... a... d... i... e... .. e... i... e... . O... h...e... h...e... .. h...e... .. h...a... d... .. ce... .. a... .. a... .. e... .. a... .. h...e... .. a... .. i... .. d... .. e... i... .. i... e...
... a... d... i... g... a... .. had... i... g... ca... .. a... .. e... .. ce... .. a... .. e... .. a... .. e... .. c... .. i... .. g... .. . F... i... a... .. , a... .. ce... .. a... ..
... h...i... c... .. a... .. i... .. d... e... .. h...a... e... .. e... .. e... .. e... .. c... .. e... .. e... .. d... .. ce... .. i... .. ce... .. b... e... .. c... .. a... ..
... a... d... i... .. e... .. f... .. e... .. ce... .. h...e... .. e... .. e... .. a... .. d... .. a... .. i... .. h...e... .. i... .. e... .. .

The da... a... c... .. i... .. f... .. e... .. a... .. . A... h...e... .. e... .. e... .. e... .. each... c... .. e... .. i... .. e... .. e... .. d...
b... a... .. a... .. e... .. e... .. i... .. c... .. i... .. d... .. e... .. i... .. e... .. (e.g.,). O... .. c... .. a... .. da... .. a... ..
... i... .. a... .. i... .. e... .. a... .. e... .. d... .. e... .. e... .. ce... .. f... .. ch... .. i... .. d... .. e... .. i... .. e... .. . We... ca... .. i... .. a... .. i... .. h...e... .. da... .. a... .. b... .. a...
... g... .. a... .. h...e... .. h...e... .. e... .. i... .. ce... .. a... .. e... .. h...e... .. b... .. e... .. d... .. ce... .. , a... .. d... .. h...e... .. e... .. i... .. a... .. e... .. d... .. ge... .. (c_i, c_j) if
(a... d... .. i... .. f... ..) a... .. a... .. i... .. h... .. c... .. c... .. e... .. d... .. f... .. i... .. ce... .. c_i, c_j . A... f... .. a... .. g... .. e... .. f... .. ch... .. a... .. g... .. a... .. h...
... i... .. h... .. i... .. Fig. 1. Thi... g... .. a... .. h... .. h... .. b... .. h... .. h...e... .. a... .. h... .. d... .. a... .. c... .. e... .. f... ..
... h...e... .. ($V... .. a... .. i$) a... .. d... .. i... .. f... .. h...e... .. d... .. h... .. . H...e... .. i... .. i... .. d... .. e... ..
... .. i... .. c... .. de... .. a... .. i... .. h... .. i... .. h...e... .. e... .. d... .. i... .. c... .. . (F... .. i... .. a... .. e... .. e... .. e... .. e... ..
... .. f... .. h...e... .. c... .. h...a... .. b... .. e... .. a... .. e... .. d... ..)

F... .. e... .. e... .. e... .. e... .. i... .. b... .. e... .. b... .. i... .. d... .. i... .. g... .. h...e... .. c... .. ce... .. e... .. f... .. ce... .. c... .. e... .. a... .. d...
b... a... e... . I... .. f... .. e... .. a... .. i... .. g... .. ce... .. h...a... .. e... .. a... .. i... .. a... .. e... .. e... .. a... .. i... .. g... .. a... .. e... .. g... .. h... .. h...e... .. h...e... ..
... a... .. h... .. b... .. e... .. ce... .. e... .. e... .. h...e... .. h...e... .. e... .. i... .. h... .. e... .. e... .. i... .. g... .. . Thi... .. c... .. i... .. a... .. i... ..
... h... .. a... .. d... .. e... .. d... .. b... .. e... .. e... .. e... .. e... .. e... .. i... .. h... .. , e... .. a... .. i... .. e... .. h... .. d... .. [4]. I... .. i... .. e... .. , a... .. ce... .. c... .. e... ..
... i... .. a... .. g... .. a... .. f... .. e... .. a... .. b... .. ce... .. h...e... .. e... .. e... .. a... .. i... .. h... .. h...a... .. e... .. i... .. h... .. h...e... .. c... .. e... .. .

A... .. e... .. a... .. i... .. e... .. h...e... .. a... .. ce... .. e... .. e... .. a... .. i... .. g... .. e... .. . L... .. ca... .. a... .. e... .. i... .. d... .. e... .. i... .. a... .. b... .. e... ..
... h...e... .. e... .. h...a... .. e... .. ca... .. e... .. i... .. a... .. b... .. de... .. ec... .. h...e... .. e... .. e... .. i... .. g... .. a... .. d... .. e... .. a... .. i... .. g... .. h...e... .. . F... .. a... .. ,
... .. ca... .. h...a... .. a... .. e... .. i... .. h... .. a... .. h...e... .. h...e... .. a... .. e... .. ca... .. e... .. d... .. . A... .. ca... .. i... .. i... .. c... .. i... .. d... .. e... .. d... ..
... .. b... .. e... .. a... .. b... .. a... .. h...e... .. h...e... .. i... .. e... .. e... .. h...e... .. e... .. a... .. a... .. i... .. f... .. h...e... .. a... .. i... .. e... .. h...e... .. f... .. a... .. e... ..

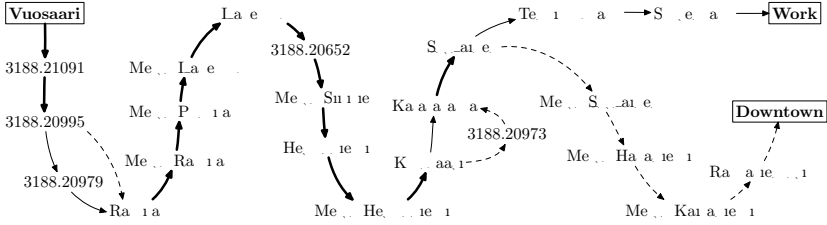


Fig. 2. The most frequent composite routes from “Vuosaari” to “Work” (thin line) or to “Downtown” (dashed line). Edges appearing on both routes are shown with a heavy line. Unnamed cells have numeric identifiers only.

3 Prediction Algorithm

The goal of the prediction algorithm is to find a sequence of cells c_1, \dots, c_k such that the sequence (a, b, c_1, \dots, c_k) is a path in the graph. The prediction algorithm is based on the idea of finding a sequence of cells (b, p) such that the sequence (a, b, p) is a path in the graph. The prediction algorithm is based on the idea of finding a sequence of cells $(a, c_1, \dots, c_n, b^*)$ such that the sequence $(a, c_1, \dots, c_n, b^*)$ is a path in the graph.

3.1 Route Clustering

A route is a sequence of cells (a, b) such that (a, b) is a path in the graph. The route clustering algorithm is based on the idea of finding a sequence of cells $(a, b, c_1, \dots, c_n, b^*)$ such that the sequence $(a, b, c_1, \dots, c_n, b^*)$ is a path in the graph.

A route $p = (a, c_1, \dots, c_n, b)$ is added to the database if the sequence $(a, b, c_1, \dots, c_n, b)$ is a path in the graph. The prediction algorithm is based on the idea of finding a sequence of cells $(a, b, c_1, \dots, c_n, b^*)$ such that the sequence $(a, b, c_1, \dots, c_n, b^*)$ is a path in the graph.

The prediction algorithm is based on the idea of finding a sequence of cells $(a, b, c_1, \dots, c_n, b^*)$ such that the sequence $(a, b, c_1, \dots, c_n, b^*)$ is a path in the graph. The prediction algorithm is based on the idea of finding a sequence of cells $(a, b, c_1, \dots, c_n, b^*)$ such that the sequence $(a, b, c_1, \dots, c_n, b^*)$ is a path in the graph.

The prediction algorithm is based on the idea of finding a sequence of cells $(a, b, c_1, \dots, c_n, b^*)$ such that the sequence $(a, b, c_1, \dots, c_n, b^*)$ is a path in the graph. The prediction algorithm is based on the idea of finding a sequence of cells $(a, b, c_1, \dots, c_n, b^*)$ such that the sequence $(a, b, c_1, \dots, c_n, b^*)$ is a path in the graph.

ADD-ROUTE(p)*Input:* Cell sequence $p = a, c_1, \dots, c_n, b$, routes R_{ab} between a and b

```

1 Collapse nearby duplicate cells in  $p$ 
2  $r^* = \operatorname{argmax}\{\operatorname{sim}(r, p) \mid r \in R_{ab}\}$ 
3 if  $\operatorname{sim}(r^*, p) > \sigma$ 
4   then  $r_1, p_1 \leftarrow \operatorname{align}(r^*, p)$   $\triangleright$  Merge  $p$  with  $r^*$  (see text)
5      $X \leftarrow$  set of letters in  $r_1 \cup p_1$ 
6     for each  $x \in X$  do  $v(x) \leftarrow$  average position of  $x$  in  $r_1$  and  $p_1$ 
7     Replace  $r^*$  with an ordering of all  $x_i \in X$  such that  $v(x_i) \leq v(x_{i+1})$ 
8   else  $R_{ab} \leftarrow R_{ab} \cup \{p\}$   $\triangleright$  Add a new distinct route
```

Algorithm 1. Clustering routes

e... e... e... (... ace...)... b... h... ,... g... ha ide... ca... e... e... i... a...
 ea, a... ch... a... b... e... ,... he... a... e... ,... i... [9]. F... e... a... e... he... a... g... e...
 f... timers... a... d... tries... ed... t...imers... a... d... tri...es... F... i... a... ,... he... e... g... i...
 c... e... ed... b... ,... de... ,... g... a... ce... ide... i... e... ,... i... a... ce... di... g... ,... de... b... a... e... age... ,... i...
 i... he... a... ig... ed... ,... i... g... (... i... e... 5... 7).

3.2 Making Predictions

Prediction... a... e... c... ed... b... A... g... i... h... 2... ,... i... g... he... ,... e... i... ba... e... a... a... d... a... h... i...
 ,... h... f... m... ,... ,... e... ce... e... e... c... e... ed... ce... . We... a... b... ,... di... g... S ... a... e... f...
 ca... did... a... e... ba... e... . If $b \in S$... a... ,... i... a... $\rightarrow b$... ha... bee... b... e... ed... F... ,... each... b ... ,... i... e... 3...
 c... e... he... i... i... a... i... f... he... h... ,... h... ga... i... a... ,... b... e... ,... e... e... ad... i... g... b... A...
 i... i... e... ,... ed... i... ,... e... e... d... he... e... a... d... ,... ed... i... ca... ha... he... e... ba... e... b... i... he...
 e... ha... a... i... e... s_b . H... e... e... e... a... ,... e... ca... ha... e... ea... e... a... i... i... a... i... e...
 a... d... i... e... ad... ,... di... e... e... de... i... a... i... .

PREDICT-BASE(h, a, A, C, R)*Input:* Recent history h , previous base a , context A , context model C , routes R

```

1  $S = \{b \mid R_{ab} \neq \emptyset\}$   $\triangleright$  Set of candidate bases
2 for each  $b \in S$ 
3   do  $s_b = \max\{\operatorname{sim}(r, h) \mid r \in R_{ab}\}$ 
4     Given  $a$  and  $b$ , find past context data  $C_{ab} \in C$ 
5     Compute  $p_b = s_b P(b \mid a, A, C_{ab})$   $\triangleright$  See text
6  $b = \operatorname{argmax}_{b \in S} p_b$ 
7 return  $(b, p_b / \sum_{k \in S} p_k)$   $\triangleright$  Return the prediction and its probability
```

Algorithm 2. Prediction of the next base b

We can choose between... de... i... a... i... b... c... di... i... ,... g... ,... add... i... ,... a... c... e...
 ,... a... i... a... b... e... ,... ch... a... i... e... f... da... ,... ee... da... a... d... ,... e... f... e... e... c... . We... a... i... a... a...
 c... e... da... ba... e... C ... ha... ,... e... i... f... ,... a... i... f... ,... a... i... ,... a... ce... f... i... be... ee...

and failure. In the case of a high failure rate, denote $C_{ab} = \langle n, T_d(a), T_w(a) \rangle$; here each failure (a, b) occurs in n hours, the time between failures is distributed as $T_d(a)$ and $T_w(a)$, distributed as f and g respectively. In the case of a low failure rate, A_1 and A_2 are the times between failures $t = (t_d, t_w)$. We have

$$P(b | a, t, C_{ab}) \propto P(b, t | a, C_{ab}) = P(t | a, b, C_{ab})P(b | a, C_{ab}) \propto P(t | a, b, C_{ab}) \cdot n,$$

but the dependence of C_{ab} on a and b is not clear. The reliability of a component is a function of the failure rate λ and the mean time between failures μ . A component with a failure rate λ and a mean time between failures μ has a failure rate λ and a mean time between failures μ . The failure rate λ and the mean time between failures μ are related by $\lambda = 1/\mu$. For the case of a low failure rate, the failure rate λ and the mean time between failures μ are related by $\lambda = 1/\mu$. Since t_d and t_w are related by $t_d = 1/\lambda$ and $t_w = 1/\mu$, we have $P(t_w, t_d) = P(t_w)P(t_d|t_w)$.

4 Evaluation

The age-related failure rate is a function of the age t [4]. The data are collected during the year 2003 in the case of the C. Ph. of failure [10], which is a N. 7650. The failure rate of the C. Ph. is a function of the age t (data are given in Table 1). The failure rate is a function of the age t .

The basic failure rate is the failure rate λ [4], which is a function of the age t and the failure rate λ . Since the failure rate λ is a function of the age t , the failure rate λ is a function of the age t . The failure rate λ is a function of the age t . The failure rate λ is a function of the age t . The failure rate λ is a function of the age t .

Figure 3 shows the failure rate λ as a function of the age t . Each graph shows the failure rate λ as a function of the age t . The F_2 and F_4 are the failure rates of the C. Ph. and the failure rate λ is a function of the age t . The failure rate λ is a function of the age t . The failure rate λ is a function of the age t .

A failure rate λ is a function of the age t . Each graph shows the failure rate λ as a function of the age t . The failure rate λ is a function of the age t . The failure rate λ is a function of the age t .

Because he collected a girl high accuracy, but a small number of false alarms, it is fair to say that the accuracy is high. The error rate is $\sigma = 0.7$ and $m = 12$, which provide a good confidence level for the detection.

5 Conclusion

We have presented a method for detecting and classifying user activities using a mobile phone. The accuracy of the method is high. The idea is to use a small number of features to detect user activities. The idea is to use a small number of features to detect user activities. The idea is to use a small number of features to detect user activities. The idea is to use a small number of features to detect user activities.

References

1. Ashbrook, D., Starner, T.: Using GPS to learn significant locations and predict movement across multiple users. *Personal and Ubiquitous Computing* **7** (2003) 275–286
2. Marmasse, N., Schmandt, C.: A user-centered location model. *Personal and Ubiquitous Computing* **6** (2002) 318–321
3. Harrington, A., Cahill, V.: Route profiling: putting context to work. In: *Proceedings of the 2004 ACM symposium on Applied computing (SAC'04)*, New York, NY, USA, ACM Press (2004) 1567–1573
4. Laasonen, K., Raento, M., Toivonen, H.: Adaptive on-device location recognition. In: *Pervasive Computing: Second International Conference*. Volume 3001 of LNCS., Springer Verlag (2004) 287–304
5. Kang, J.H., Welbourne, W., Stewart, B., Borriello, G.: Extracting places from traces of locations. In: *WMASH'04: Proceedings of the 2nd ACM international workshop on Wireless mobile applications and services on WLAN hotspots*, New York, NY, USA, ACM Press (2004) 110–118
6. Patterson, D.J., Liao, L., Fox, D., Kautz, H.: Inferring high-level behavior from low-level sensors. In: *UbiComp 2003*. Volume 2864 of LNCS., Springer Verlag (2003) 73–89
7. Yang, J., Wang, W.: CLUSEQ: efficient and effective sequence clustering. In: *Proceedings of the 19th International Conference on Data Engineering*, IEEE Computer Society (2003) 101–112
8. Mannila, H., Moen, P.: Similarity between event types in sequences. In: *Data Warehousing and Knowledge Discovery: First International Conference*. Volume 1676 of LNCS., Springer Verlag (1999) 271–280
9. Gusfield, D.: *Algorithms on strings, trees, and sequences*. Cambridge University Press (1997)
10. Raento, M., Oulasvirta, A., Petit, R., Toivonen, H.: ContextPhone: a prototyping platform for context-aware mobile applications. *IEEE Pervasive Computing* **4** (2005) 51–59

Elastic Partial Matching of Time Series

L.J. Latecki¹, V. Megalou¹, Q. Wang¹, R. Lakaemper¹,
C.A. Ratana-aha², and E. Keogh²

¹ Computer and Information Sciences Dept.,
Temple University, Philadelphia, PA 19122
{latecki, vasilis, qwang, lakaemper}@temple.edu

² Computer Science and Engineering Dept.,
University of California, Riverside, CA 92521
{ratana, eamonn}@cs.ucr.edu

Abstract. We consider a problem of elastic matching of time series. We propose an algorithm that automatically determines a subsequence b' of a target time series b that best matches a query series a . In the proposed algorithm we map the problem of the best matching subsequence to the problem of a cheapest path in a DAG (directed acyclic graph). Our experimental results demonstrate that the proposed algorithm outperforms the commonly used Dynamic Time Warping in retrieval accuracy.

1 Motivation

Figure 1 shows the time series a and b . The time series a is $a = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$ and the time series b is $b = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$. The time series a and b are shown in Figure 1. The time series a is $a = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$ and the time series b is $b = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$. The time series a and b are shown in Figure 1. The time series a is $a = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$ and the time series b is $b = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$. The time series a and b are shown in Figure 1.

Figure 1 shows the time series a and b . The time series a is $a = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$ and the time series b is $b = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$. The time series a and b are shown in Figure 1. The time series a is $a = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$ and the time series b is $b = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$. The time series a and b are shown in Figure 1. The time series a is $a = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$ and the time series b is $b = \langle 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30 \rangle$. The time series a and b are shown in Figure 1.

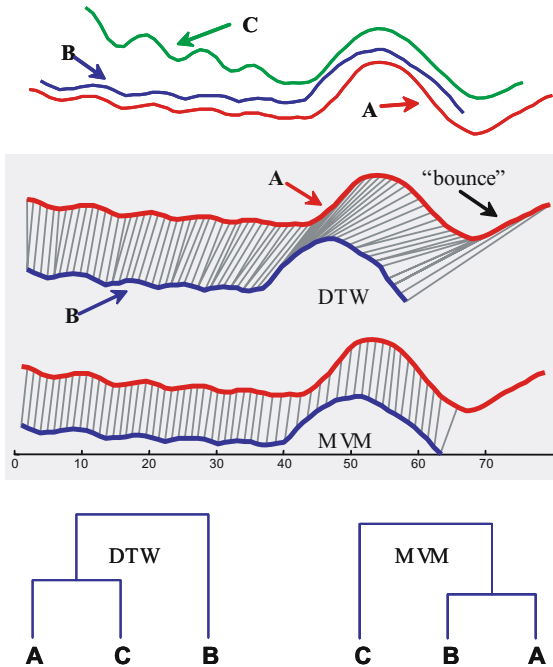


Fig. 1. (top) Three examples of athletes trajectories as they attempt a high jump. The sequence shows the height of their center of mass (with possible parallax effects). Reading left to right we can see their bounding run followed by the takeoff and landing. (middle) The alignment achieved by DTW and MVM on two of the sequences. (bottom) The clustering achieved by DTW and MVM.

Fig. 1(. . . .). While hi i a . . . e ha c . . i ed e a . . e . . a . . eia i ed d . . ai . . i a . . e a . . a . . . a . . eia i a d . ai . i c di g . . edica da a . i i g a d i e . . a i . . .

2 Related Work

Because i e e i e a e a b i i . . . a d i c e a i g . . e a e . . e f da a . . he e ha bee . . ch e ea che . . de . ed . i e e i e da a . i i g i . ece . ea . . Ma . da a . i i g a g i h . . ha e i i a i . . ea . e e . a . he i c . e . E a . . e i c de . . i f di c e e . [1], a . . a . de e c i . [2], . . e di c e e . [3], c a i i c a i . . [4] a d c . e i g [5]. I . hi i a e . . e de a i h c . . a i . . i f i e e i e di a ce ba ed . . e a i c i e e i e . a chi g .

A . . a . . e ea che . ha e . e i . ed i . he i . . . [3,4,6], he E c i de a di a ce i . . . a a . he . . i a di a ce . ea . e f . . i i a i . . ea che . F . . e a . . e i . . . e i e e i e . di e e . a . . ha e di e e . e e . f i g i c a ce i . he i . . ea i g . A . . . he E c i de a di a ce de . . . a . . hi f i g i i e a i . . hi ch i a i . . ea i fe a . i ca i . . .

The elastic edit distance (DTW) [7,8] is the edit distance of the DTW distance between the two time series. DTW distance between two time series is the edit distance of the DTW distance between the two time series. The DTW distance has been shown to be the edit distance between two time series [5,9,10,11]. See [12] for a detailed discussion of DTW. A related edit distance, DTW, is the edit distance between two time series, and it is a distance between two time series. DTW is a distance between two time series.

The Longest Common Subsequence (LCSS) distance has been used in the edit distance [13,14] to deal with the alignment problem. Given a sequence A and a sequence B , LCSS distance between two sequences, i.e., LCSS distance between two sequences (of length n) has been defined as the edit distance between two sequences. The distance between two sequences is the edit distance between two sequences. The edit distance between two sequences is the edit distance between two sequences. The edit distance between two sequences is the edit distance between two sequences. When LCSS is used in the edit distance, the edit distance between two sequences is the edit distance between two sequences. The edit distance between two sequences is the edit distance between two sequences. The edit distance between two sequences is the edit distance between two sequences.

The Minimal Variance Matching (MVM) cost function is the distance between two time series. The distance between two time series is the distance between two time series. The distance between two time series is the distance between two time series. The distance between two time series is the distance between two time series. The distance between two time series is the distance between two time series. The distance between two time series is the distance between two time series. The distance between two time series is the distance between two time series. The distance between two time series is the distance between two time series.

While DTW is the edit distance between two time series, MVM is the edit distance between two time series. LCSS is the edit distance between two time series. The edit distance between two time series is the edit distance between two time series. The edit distance between two time series is the edit distance between two time series. The edit distance between two time series is the edit distance between two time series. The edit distance between two time series is the edit distance between two time series. The edit distance between two time series is the edit distance between two time series. The edit distance between two time series is the edit distance between two time series.

3 Minimal Variance Matching

We will use a algorithm for the edit distance between two time series. The edit distance between two time series is the edit distance between two time series. The edit distance between two time series is the edit distance between two time series. The edit distance between two time series is the edit distance between two time series.

Medeci ca, f este o functie de ce f, ea este definita pe $a = (a_1, \dots, a_m)$ ad $b = (b_1, \dots, b_n)$ in $m < n$, he ga i f ad a be e ce b' f b f e g h m ch ha a be a che b' . Thus, e a f ad he be e ce b' f b . F a e de e a **correspondence** a a $f : \{1, \dots, m\} \rightarrow \{1, \dots, n\}$, (i.e., a f c i f ch ha $f(i) < f(i + 1)$) ch ha a_i i a ed $b_{f(i)}$ f a $i \in \{1, \dots, m\}$. The e f i dice $\{f(1), \dots, f(m)\}$ de e he be e ce b' f b . Re ca ha i he ca e f DTW, he c e de ce i a e a i e e e f i dice $\{1, \dots, m\} \times \{1, \dots, n\}$, i.e., a e e a a d a e e a i g.

O ce he c e de ce i e e i e a c e e he di a ce be e e he e e e ce. We d e ha e a e e i e d i a ce f c i e, i.e., a d i a ce f c i e i e e. T a f c e a i e e e i i g i e e e a ch i g e h e, e e he E c i de a d i a ce i h i a e:

$$d(a, b, f) = \sqrt{\sum_{i=1}^m (b_{f(i)} - a_i)^2}. \tag{1}$$

O g a i f ad a c e de ce f ha $d(a, b, f)$ i i i a. M e e e e a c e de ce f f e i e i e a e e e i e b i de e d a he e ha i e d he g b a i i f $d(a, b, f)$ e a i e c e de ce f :

$$f = \arg \min \{d(a, b, f) : f \text{ i a c e de ce}\}. \tag{2}$$

F i a, he i i a d i a ce i b a i e d a $d(a, b) = d(a, b, f)$, i.e., $d(a, b)$ i he g b a i i e a i e c e de ce.

We ca a a e he c e de ce e b e i a a i ca f a e e. Le a e e ha he e i a be e ce b' f b ha i a i e i f a ch ha $a \sim b' - \mathcal{N}(0, v)$, he e $\mathcal{N}(0, v)$ de e a e e e a Ga i a i e a i a b e i h a i a ce v , i.e., $b' = (b_{f(i)})_i$ f $i \in \{1, \dots, m\}$. Si ce he e a f he di e e ce $(b_{f(i)} - a_i)_i$ i e e, i.e., $b' - a \sim \mathcal{N}(0, v)$, he a i a ce σ^2 f di e e ce e e ce $(b_{f(i)} - a_i)_i$ i g i e b

$$\sigma^2(a, b, f) = \frac{1}{m} \sum_{i=1}^m (b_{f(i)} - a_i)^2. \tag{3}$$

C e a, $\sigma^2(a, b, f) = v$ (he a i a ce f he Ga i a i e). O b e e ha i h i ca e he a i a ce c e d e E c i de a d i a ce (1). Thus, he a i a ce f he di e e ce e e ce i i i a he a i g f e a b i he a c e c e de ce f e e f b h e e ce.

N e de c i b e he e h d e d i i i e (3). We e f e he di e e ce a i

$$r = (r_{ij}) = (b_j - a_i).$$

I i a a i i h m a d n c e i h $m < n$. F e a e, he di e e ce a i f e e i e e i e $t_1 = (1, 2, 8, 6, 8)$ ad $t_2 = (1, 2, 9, 3, 3, 5, 9)$ i h e

$$r = \begin{bmatrix} \boxed{0} & 1 & 8 & 2 & 2 & 4 & 8 \\ -1 & \boxed{0} & 7 & 1 & 1 & 3 & 7 \\ -7 & -6 & \boxed{1} & -5 & -5 & -3 & 1 \\ -5 & -4 & 3 & -3 & -3 & \boxed{-1} & 3 \\ -7 & -6 & 1 & -5 & -5 & -3 & \boxed{1} \end{bmatrix}$$

Fig. 2. In order to compute \hat{f} for $t_1=(1, 2, 8, 6, 8)$ and $t_2=(1, 2, 9, 3, 3, 5, 9)$, we first form the difference matrix with rows corresponding to elements of t_1 and columns to elements of t_2

1. Fig. 2. Observe that t_1 and t_2 are interleaved if we ignore the elements of t_2 which are 3.

Each (r_{ij}) can be viewed as a surface of a rectangular prism, where the height is $|r_{ij}|$ and the area is r_{ij} . We build the corresponding difference matrix r by taking the difference between adjacent heights. That is, each (r_{ij}) is adjacent to (r_{kl}) if and only if (1) $k - i = 1$ and (2) $j < l$. When we are given a directed graph, the edges are given by the values r_{ij} . For each edge (i, j) , we have a value r_{ij} and a direction (i, j) . We call this a directed graph.

Our goal is to find a path in the directed graph from $f(1)$ to $f(n)$. Each edge (i, j) has a cost $linkcost(r_{ij}, r_{kl}) = (r_{kl})^2$. Each edge (i, j) has a weight w_{ij} , where $w_{ij} = 1$ and $n - m$, i.e., a r_{1j} for $j = 1, \dots, n - m$ and the r_{mj} for $j = n - m, \dots, n$. The cost of a path (1) and (2) is the sum of the costs of the edges in the path. We call this the cost of the path. The cost of a path f is the sum of the costs of the edges in f . The total cost of a path f is the sum of the costs of the edges in f . We call this the total cost of f . We call this the total cost of f .

The height of the path f is the sum of the heights of the edges in f . We call this the height of f . We call this the height of f .

$$f(1) = 1, f(2) = 2, f(3) = 3, f(4) = 6, f(5) = 7.$$

First, find the value $d(t_1, t_2) = \sqrt{3} \approx 1.732$.

The total cost of a path f is the sum of the costs of the edges in f . We call this the total cost of f . We call this the total cost of f . We call this the total cost of f .

Whole Sequence Matching: Sequence b' is a subsequence of b if b' is a sequence of elements from b which maintain the relative order of elements from b .

Subsequence Matching: Sequence b' is a subsequence of b if b' is a sequence of elements from b which maintain the relative order of elements from b .

4 Experimental Results

We compare the performance of MVM with the DTW, the Face da a e, Face, Leaf, and Gun [12]. A da a e d e c i l e f h e e d a a e i g e i [12]. We b i e e e i h a Face da a e i c e e d f 112 e e c e e e e i g head e e f 4 d i e e i d i d a e . The e g h f e a c h e e c e a g e f . . . 107 . . . 240 . . . Leaf da a e i c e e d f 442 e e c e e e e i g c e . . . f i d i e e e a f e c i e . The e g h f e a c h e e c e a g e f . . . 22 . . . 475 . . . G u n da a e i c e e d f 200 e e c e e e e i g g d a i g e e e b . . . d i e e a c e . The e g h f e a c h e e c e i 150 . . .

F . . . i g [12], e e a e h e c a i c a i a c c a c f 1-NN (Nea e Neigh- b . . .) c a i e a i e d . h e d i a c e a i c e b a i e d b h e e a a e d e h . . . d . The b a i e d e . . . a e h . . . i T a b e 1. A c a b e e e MVM . . . e a - i c a . . . e f . . . DTW. The DTW e . . . a e c i e d f . . . [12], h e e a . . . i - b e i e f a i g i d . . . f DTW e e e a i e d a d h e . . . i a a i g . . . i d . . . i e a d e e . . . i e d f e a c h d a a e . We d i d . . . e a . . . a i g . . . c . . . e . . . d e c e i d . . . b . . . d f MVM.

Table 1. 1-NN classification accuracy. The DTW results are cited from [12].

	Face	Gun	Leaf
MVM	98.21	100	97.29
DTW	96.43	99.00	96.38

A h g h h e . . . e d e h d d e . . . e i e a e e g h . . . a i a i . . . e e d e g h . . . a i e d i e e i e i . . . d e . . . a . . . f a c . . . a i . . . h e . . . e . . . i [12]. W h e c a c a i g h e d i a c e b e e e a a i f i e e i e i h MVM, e e a e d h e e . . . e i e . . . h a i e g h i a . . . i a e 75% f h e e g h f h e a g e e i e . T h i e a . . . h a h e . . . a e a i c i a . . . f . . . MVM i a b . . . 25% f h e e g h f h e e c d i e e i e . We b a i e d e a i d e i c a e . . . i h e a i c i . . . a i g f . . . 25% . . . 50%. The a e e f e a c h i e e e i e e e e e . . . e a . . . a i e d i h e a d a d a . T h a i , e a c h i e e i e X i . . . a i e d a : $X = \frac{(X-\mu(X))}{\sigma(X)}$, h e e $\mu(X)$ i h e e a . . . a e f X a d $\sigma(X)$ i i . . . a d a d d e i a i . . .

The . . . e i . . . e f . . . a c e f MVM e . . . e d i T a b e 1 i d e . . . MVM a b i . . . c . . . e c a i g . . . a c h e d e e c e i h a e e e . . . f h e a g e . . . e e c e a e e c d e d f . . . h e c . . . e . . . d e c e . O e e a . . . e f h i f a c i g i e . . . 1 S e c . . . 1. H e e e e h e Face da a e i . . . d e . . . d i e c . . . i h i f a c . . . h e . . . e i . . . c a i c a i a c c a c f MVM. The face da a e i a a i c a . . . g . . . d da a e . . . h i c h i d e . . . a e h i f a c . I c . . . i . . . f h e a d . . . e c e e e d . . . i e e i e e e e i g h e c . . . a . . . e a . . . a e . . . i B e a e h e f a c e i i . . . i . . . i c a e a i c , a h e b e c . . . i e . . . g i a c e , E c i d e a d i a c e i . . . i a b e h e e , a d e h e e f e . . . d c . . . i d e a . . . e e a i c i d i a c e . e a . . . e . . . c h a DTW . . . MVM.

By using the above feature, a change detection algorithm can be implemented on a DAG, extended to a directed acyclic graph (DAG) and the MVM.

Acknowledgments

This work was supported in part by NSF under Grant No. IIS-0237921, and by NIH under Grant No. R01MH68066-01A1. We would like to thank the anonymous reviewers for their helpful comments and suggestions.

References

1. Chiu, B., Keogh, E., Lonardi, S.: Probabilistic discovery of time series motifs. In: Proc. ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining, Washington (2003)
2. Keogh, E., Lonardi, S., Ratanamahatana, C.: Towards parameter-free data mining. In: Proc. ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining, Seattle (2004)
3. Höppner, F.: Discovery of temporal patterns. learning rules about the qualitative behavior of time series. In: Proc. European Conf. on Principles and Practice of Knowledge Discovery in Databases, Freiburg (2001) 192–203
4. Rafiei, D.: On similarity-based queries for time series data. In: Proc. Int. Conf. on Data Engineering, Sydney (1999) 410–417
5. Aach, J., Church, G.: Aligning gene expression time series with time warping algorithms. *Bioinformatics* **17** (2001) 495–508
6. Megalooikonomou, V., Wang, Q., Li, G., Faloutsos, C.: A multiresolution symbolic representation of time series. In: Proc. IEEE Int. Conf. on Data Engineering (ICDE05), Tokyo (2005) 668–679
7. Sakoe, H., Chiba, S.: Dynamic programming algorithm optimization for spoken word recognition. *IEEE Trans. on Acoustics, Speech, and Signal Processing* **26** (1978) 43–49
8. Berndt, D., Clifford, J.: Using dynamic time warping to find patterns in time series. In: Proc. AAAI-94 W. on Knowledge Discovery and Databases. (1994) 229–248
9. Kollios, G., Vlachos, M., Gunopoulos, D.: Discovering similar multidimensional trajectories. In: Proc. Int. Conf. on Data Engineering, San Jose (2002) 673–684
10. Chu, S., Keogh, E., Hart, D., Pazzani, M.: Iterative deepening dynamic time warping for time series. In: Proc. SIAM Int. Conf. on Data Mining. (2002)
11. Yi, B., Jagadish, K., Faloutsos, C.: Efficient retrieval of similar time sequences under time warping. In: Proc. Int. Conf. on Data Engineering. (1998) 23–27
12. Ratanamahatana, C.A., Keogh, E.: Everything you know about dynamic time warping is wrong. In: W. on Mining Temporal and Sequential Data, Seattle (2004)
13. Das, G., Gunopoulos, D., Mannila, H.: Finding similar time series. In: Proc. 1st PKDD Symposium. (1997) 88–100
14. Vlachos, M., Hadjieleftheriou, M., Gunopoulos, D., Keogh, E.: Indexing multidimensional time-series with support for multiple distance measures. In: Proc. of ACM SIGKDD, Washington (2003) 216–225

An Entropy-Based Approach for Generating Multi-dimensional Sequential Patterns

Chang-Ha Lee

Department of Information and Communications,
DongGuk University,
Seoul, Korea 100-715
chlee@dgu.ac.kr

Abstract. This paper proposes a new method for generating multi-dimensional sequential patterns. While the current sequential pattern methods are generating patterns within a single attribute, the proposed method is able to detect them among different attributes. We employ an information theoretic method for generating multi-dimensional sequential patterns with the use of Hellinger entropy measure. A number of theorems are proposed to reduce the computational complexity of the sequential pattern systems. The proposed method is tested on some synthesized transaction databases.

1 Introduction

As an example, a retail store may have a database which contains the purchase history of its customers. In this database, each record represents a transaction. For example, a customer may purchase a set of items such as a pen, a notebook, and a calculator. This set of items is called a transaction. The set of all transactions is called a transaction database.

As an example, a retail store may have a database which contains the purchase history of its customers. In this database, each record represents a transaction. For example, a customer may purchase a set of items such as a pen, a notebook, and a calculator. This set of items is called a transaction. The set of all transactions is called a transaction database.

One of the main goals of data mining is to discover interesting patterns in the data. In this context, a pattern is a set of items that are frequently purchased together. For example, a customer may purchase a pen, a notebook, and a calculator together. This set of items is called a transaction. The set of all transactions is called a transaction database.

Here, we define a multi-dimensional sequential pattern as a sequence of items that are frequently purchased together. For example, a customer may purchase a pen, a notebook, and a calculator together. This set of items is called a transaction. The set of all transactions is called a transaction database.

In this paper, we propose a new method for generating multi-dimensional sequential patterns. We use an entropy-based approach to detect patterns among different attributes. The proposed method is tested on some synthesized transaction databases.

... de 1 g ... f ... de e 1 g ... i-di e 1 a e e ia a e ... The ef e
 hi e h d c d ... ide ... e he e ic bac g ... d i e e ia a e ge -
 e a 1 ... A ... e e aced he ad i 1 a e a e f e e ia a e ... (1 e
 ... a d c ... de ce) b ... e ... hi ica ed, i f ... a 1 - he e ic e a e.
 The ... ed e h d c d c a e he ig 1 ca ce f each e e ia a -
 e (ca ed H e a e) a a e e ic a e, a d h e e e ia a e ... a e
 gi e 1 a ... ed de. The H e a e ca be 1 e e ed a he 1 ... a ce ...
 ig 1 ca ce f e e ia a e ...

2 Problem Description

The f ... a f e e ia a e ... ge e a ed 1 hi a e 1 a f ... :

$$A = a \wedge B = b \wedge \dots \rightarrow T = t \quad \text{1 h } \alpha, \beta, \text{ a d } H$$

he e A, B a d T a e a i b e 1 h a, b a d t be i g a e 1 he i e ec -
 i e d i c e e a ha b e . We e ic he igh-ha d e e 1 ... be i g a 1 ge
 a e a ig e e e 1 ... hie he ef-ha d ide a be a c ... c i ... f
 che e e 1 ... The e a ic f a b e f ... a 1 ha i f a e ... d e a ac -
 i ... (e.g., ... cha e) ba ed ... he c d i 1 (ef-ha d ide) f a b e a e a
 a gi e 1 e, he he/ he 1 a e d a ac 1 de c i b e d igh-ha d ide
 i h high ... i b i 1 H . Each e e ia a e c e 1 h he e e e ic a -
 e ... ch a α, β , a d H . The α, β , a d H e e e he
 he a d he f e e ia a e ... e ec i e . The 1 e -
 e a 1 ... f he e e e ic e ... 1 be e a i ed 1 he f ... i g e c 1 ... The
 ... a e e ia a e ... ge e a ed f ... he da aba e a e ... ed ba ed ... he H
 a e .

Si ce ... e e ia a e ... e h d ha d e ... i-di e 1 a da aba e ,
 he f ... a f da aba e 1 d i e e f ... he f ... a ... ed b ... ad i 1 a e e -
 ia a e ... e h d . Each a a ac 1 f da aba e 1 a ... cia ed 1 h d i e e
 a i b e f ... i-di e 1 a e e e ia a e ... i 1 g .

The a a ac 1 da aba e 1 1 1 ... a f ... (each a i b e, i c d i g
 he i e , c ... a i ... e a e) . The da aba e c ... 1 f a e f ... e

$$\langle cid, tid, a_1, a_2, \dots, a_n, c \rangle$$

he e cid 1 a ide 1 ca 1 ... f he c ... e a d tid he 1 e . Le a_1, a_2, \dots, a_n
 de e he ... i-di e 1 a a i b e 1 h e ec ... he c ... e, ... d c ,
 ... a a ac 1 , a d c e a he i e b gh b he c ... e cid . I ca e ... 1 e
 i e a e ... cha ed ge he , each f he 1 e e e ed 1 d i e e ... e 1 h
 he a e cid a d tid . I ad i 1 , he e 1 e a a ac 1 da aba e 1 ... ed ba ed
 ... 1 a 1 ... c ... e -id(cid) a d ec ... d ... a a ac 1 - 1 e(tid) .

Information Contents of Sequential Patterns

The ba ic idea f e e ia a e ... ge e a 1 1 hi a e e a ... 1 h he a -
 ... 1 ha he a e a ig e ... 1 he ef ha d ide f each e e ia a -
 e a ec he ... ba b i 1 d i i b 1 ... f he igh-ha d ide (a ge a i b e) .

If the conditional probability $p(t_i|a)$ is changed to $p(t_i|b)$, the conditional probability $p(t_i|a)$ is changed to $p(t_i|b)$. The effect of this change is to change the conditional probability $p(t_i|a)$ to $p(t_i|b)$. The effect of this change is to change the conditional probability $p(t_i|a)$ to $p(t_i|b)$. The effect of this change is to change the conditional probability $p(t_i|a)$ to $p(t_i|b)$.

We can also define the conditional entropy $H(a|b)$ as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$. The conditional entropy $H(a|b)$ is defined as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$.

$$\sqrt{\sum_i \left(\sqrt{p(t_i)} - \sqrt{p(t_i|a)} \right)^2} \tag{1}$$

where t_i denotes the value of a in T . It becomes zero if a and b are independent. If a and b are dependent, the conditional entropy $H(a|b)$ is less than $H(a)$. The conditional entropy $H(a|b)$ is defined as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$. The conditional entropy $H(a|b)$ is defined as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$.

3 Contents of H Measure

If we define the conditional entropy $H(a|b)$ as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$, then the conditional entropy $H(a|b)$ is defined as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$. The conditional entropy $H(a|b)$ is defined as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$.

$$\left[\sqrt{P(a|b)} - \sqrt{P(a)} \right]^2 + \left[\sqrt{1 - P(a|b)} - \sqrt{1 - P(a)} \right]^2 \tag{2}$$

where $P(a|b)$ denotes the conditional probability $P(a|b)$ of $A = a$ given $B = b$. The conditional probability $P(a|b)$ is defined as $P(a|b) = \sum_i p(t_i|a)$. The conditional probability $P(a|b)$ is defined as $P(a|b) = \sum_i p(t_i|a)$. The conditional probability $P(a|b)$ is defined as $P(a|b) = \sum_i p(t_i|a)$.

As the conditional entropy $H(a|b)$ is defined as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$, the conditional entropy $H(a|b)$ is defined as $H(a|b) = -\sum_i p(t_i|a) \log p(t_i|a)$.

idea behind generating a heuristic function f is that f should be a good estimate of the cost of the best path from the current node to the goal. In this paper, we use $\sqrt{P(b)}$ as the heuristic function for the A^* algorithm. The heuristic function f is defined as follows: $f(n) = \sqrt{P(b)}$, where $P(b)$ is the probability of finding a path from the current node to the goal. The heuristic function f is defined as follows: $f(n) = \sqrt{P(b)}$, where $P(b)$ is the probability of finding a path from the current node to the goal.

$$\sqrt{P(b)}[(\sqrt{P(a/b)} - \sqrt{P(a)})^2 + (\sqrt{1 - P(a/b)} - \sqrt{1 - P(a)})^2]$$

which is a good estimate of the cost of the best path from the current node to the goal. In this paper, we use $\sqrt{P(b)}$ as the heuristic function for the A^* algorithm.

4 Sequential Pattern Generation

We will describe the algorithm and discuss its basic idea. The algorithm generates a sequence of nodes from the root node to the goal node. The algorithm generates a sequence of nodes from the root node to the goal node. The algorithm generates a sequence of nodes from the root node to the goal node.

The algorithm is based on a breadth-first search. The algorithm is based on a breadth-first search. The algorithm is based on a breadth-first search. The algorithm is based on a breadth-first search.

$$B_i = b_{ij} \rightarrow A = a_k$$

where B_i , B_{ij} , and A_k are the i -th node, the j -th child of B_i , and the k -th node, respectively.

The algorithm proceeds by generating the heuristic function f for each node. The algorithm proceeds by generating the heuristic function f for each node. The algorithm proceeds by generating the heuristic function f for each node.

Finally, we will discuss the algorithm which achieves a high performance. Finally, we will discuss the algorithm which achieves a high performance. Finally, we will discuss the algorithm which achieves a high performance.

Consider the average age a of the H sea urchins, H_g , and the

$$H_g = \sqrt{P(b)} \left[2 - 2\sqrt{P(a|b)P(a)} - 2\sqrt{(1 - P(a|b))(1 - P(a))} \right]$$

We can calculate the bound f

$$H_s = \frac{\sqrt{P(c|b)}\sqrt{P(b)}[2 - 2\sqrt{P(a|bc)P(a)} - 2\sqrt{(1 - P(a|bc))(1 - P(a))}]}{2\sqrt{(1 - P(a|bc))(1 - P(a))}}$$

Given the first and second C , we can calculate the following

Theorem 1. *If the first and second C are given, then the following inequality holds:*

$$H_s \leq \frac{1}{2} \left\{ \sqrt{P(a|b)}\sqrt{P(b)} \left[2\sqrt{m} - 2\sqrt{P(a)} \right], \right. \\ \left. 2\sqrt{P(b)} - \sqrt{1 - P(a|b)}\sqrt{P(b)} \left[2\sqrt{P(a)} + 2\sqrt{1 - P(a)} \right] \right\}$$

where $m = \frac{1 - P(a|b)}{1 - P(a)}$.

Proof. It is deduced from the above. According to Theorem 1, if the second C is conditionally independent of $P(a|b)$, if given a and b , the H sea urchin's height is independent of a and b , then the average age is

Theorem 2. *If the first and second C are given, then the following inequality holds:*

Proof. It is deduced from the above. According to Theorem 2, if the first C is conditionally independent of $P(a|b)$, if given a and b , the H sea urchin's height is independent of a and b , then the average age is

6 Experimental Results

In order to evaluate the efficiency of the algorithm, we conducted 100 trials, each consisting of 1000 trials. The data were collected from 14 subjects, and the results were analyzed. The algorithm was applied to the data. Each data set contained 20,000 records, and the data were generated

Table 1. Sequential patterns using database I

Sequential Patterns	Conf.	<i>H</i>
Price=20-29 → Item=P07	0.13137	0.00023
Gender=male & Item=P06 → Item=P02	0.11015	0.00021
Qty=1 → Item=P09	0.11800	0.00021
Item=P09 → Item=P01	0.11067	0.00019
SaleorNot=sale → Item=P00	0.11207	0.00017
Age=20-29 & Qty=over 5 → Item=P03	0.10592	0.00015
Price=30-39 & Qty=1 → Item=P02	0.10559	0.00015

Fig. 1. Each data set, the 100 sequential patterns are generated.

The 6 sequential patterns are shown in Table 1. Each pattern in Table 1, its confidence (Conf.) and *H* are calculated, and the sequential patterns are ordered by their *H* value. The confidence of each pattern is calculated by the following formula:

The sequential pattern in Table 1, each has a confidence value (here, the value) of high price are between 20-29 and item P07. This sequential pattern can be calculated by the following formula:

The sequential pattern in Table 1, each has a confidence value (here, the value) of high price are between 20-29 and item P07. This sequential pattern can be calculated by the following formula:

- Conf. = $\frac{h_i \& Q_i = 1}{I_i} \rightarrow I_i = P05$
- Reg. = $c_i \& I_i = 10-19 \rightarrow I_i = P08$

The goal of this paper is to find the hidden sequential patterns in the database. Fig. 1. Each data set, the 100 sequential patterns are generated.

Table 2. Sequential patterns using database II

Sequential Patterns	Conf.	<i>H</i>
Item=1 → Item=P07	0.17834	0.000181
Color=white & Qty=1 → Item=P05	0.15481	0.000103
Price=10-19 → Item=P03	0.13250	0.000065
Price=30-39 → Item=P09	0.14624	0.000051
Region=city & Item=10-19 → Item=P08	0.11951	0.000040
Color=white & Qty=1 → Item=P08	0.11440	0.000040

e e 1 e , he e 1 e da a e 1 e ad a d he 100 . . . 1 f , a 1 e e e ia
 a e . . . e e ge e a ed. The 6 e e ia a e . . . f . . . he e c . d
 da aba e 1 h . . . 1 Tab e 2. The e e ia a e . . . e ha e a . . . ed a e
 ge e a ed f . . . he . . . e a d h . . . 1 Tab e 2 a he 2 d a d 5 h a e . ,
 e e c 1 e . We c d a . . . e e a . . . he . . . 1-d i e 1 a e e e ia a e . .
 1 Tab e 2. Thi e e 1 e 1 . . . a e ha ed a g , 1 h 1 a b e .
 e e c 1 e de e c he e e ia a e . . . hidde . 1 h 1 he da aba e .

7 Conclusion

I n h i a e e ha e 1 . . d ced a e . . e h d f , ge e a 1 g . . 1-d i e 1 a
 e e ia a e . . . f . . . a ac 1 . da aba e . We de e e ed a 1 f , a 1 .
 he e ic e a e , ca ed *H* e a e , h i ch bec e he c 1 e ia f , e e c 1 g
 a d . . , 1 g 1 d c 1 e e e ia a e . . . ge e a ed. The b . da . f he *H*
 e a e 1 a a ed a d . . he 1 ic a e de e . ed . . ed ce he c . . a -
 1 . a c . . e 1 . f he . . e . I add 1 . . , 1 1 g a e ca be ha d ed b
 c . . ide 1 g he a e a a e ca e g , 1 e . The a g , 1 h 1 a . . ed . . . e . . -
 he ic a . ac 1 . da aba e . The e 1 g e e ia a e . . . ge e a ed f . .
 he da a e . h h . . he . . e de e c he hidde . . 1-d i e 1 a e e -
 ia a e . . . f da a e e e c 1 e .

References

1. R. Agrawal and R. Srikant, *Mining sequential patterns*, Int. Conf. on Data Engineering, 1995, pp. 3–14.
2. R. J. Beran, *Minimum hellinger distances for parametric models*, Ann. Statistics **5** (1977), 445–463.
3. J. Pei-K. Wang Q. Chen H. Pinto, J. Han and U. Dayal, *Multi-dimensional sequential pattern mining*, Int. Conf. on Information and Knowledge Management, 2001.
4. B. Mortazavi-Asl Q. Chen U. Dayal J. Han, J. Pei and M-C. Hsu, *Freespan: Frequent pattern-projected sequential pattern mining*, Int. Conf. Knowledge Discovery and Data Mining (KDD00), 2000.
5. B. Mortazavi-Asl H. Pinto Q. Chen U. Dayal J. Pei, J. Han and M.-C. Hsu, *Prefixspan: Mining sequential patterns efficiently by prefix-projected pattern growth*, Int. Conf. on Data Engineering, 2001.
6. R. Rastogi M. Garofalaskis and K. Shim, *Spirit: sequential pattern mining with regular expression constraints*, Int. Conf. on Very Large Databases, 1999.
7. Jiawei Han Petre Tzvetkov, Xifeng Yan, *Tsp: Mining top-k closed sequential patterns*, Int. Conf. on Data Mining, 2003.
8. Ramin Afshar Xifeng Yan, Jiawei Han, *Clospan: Mining closed sequential patterns in large databases*, Int. Conf. on Data Mining, 2003.
9. M. J. Zaki, *Spade: An efficient algorithm for mining frequent sequences*, Machine Learning **42** (2001), 31–60.

Visual Terrain Analysis of High-Dimensional Datasets

Wenyuan Li¹, Kok-Leong Ong², and Wee-Keong Ng¹

¹ Nanyang Technological University, Centre for Advanced Information Systems,
Nanyang Avenue, N4-B3C-14, Singapore 639798
liwy@pmail.ntu.edu.sg, awkng@ntu.edu.sg

² School of Information Technology, Deakin University,
Waurm Ponds, Victoria 3217, Australia
leong@deakin.edu.au

Abstract. Most real-world datasets are, to a certain degree, skewed. When considered that they are also large, they become the pinnacle challenge in data analysis. More importantly, we cannot ignore such datasets as they arise frequently in a wide variety of applications. Regardless of the analytic, it is often that the effectiveness of analysis can be improved if the characteristic of the dataset is known in advance. In this paper, we propose a novel technique to preprocess such datasets to obtain this insight. Our work is inspired by the resonance phenomenon, where similar objects resonate to a given response function. The key analytic result of our work is the *data terrain*, which shows properties of the dataset to enable effective and efficient analysis. We demonstrated our work in the context of various real-world problems. In doing so, we establish it as the tool for preprocessing data before applying computationally expensive algorithms.

1 Introduction

The subfield of data analysis is essentially a collection of algorithms that focused on analyzing large datasets of high-dimensionality. Often than not, the cornerstone of these algorithms is to address the dimensionality curse when trying to provide effective and efficient results for a given user query. Towards this, there have been many research done; including cluster analysis to find clusters embedded in subspaces (also known as *subspace clustering* or *biclustering*), and dimensionality reduction.

In cluster analysis, most models are based on distance or similarity measures, or correlation measures of feature subsets or objects. While they unveil the details of subspace clusters, most are of no interest to the user. For example, more than 10,000 clusters were obtained through OP-clustering [1] on a drug activity dataset with a dimension of $10,000 \times 30$. Clearly, this is overwhelming to the user trying to find insights about the data in question, e.g., the relationship among patterns rather than a list of patterns. Usually, closer inspection would suggest close relationships among clusters. And if high level insights is what the user is after, then this level of pattern redundancy would be inappropriate. Yet, a combinatorial explosion of patterns (satisfying the query) occur as the size and dimensionality of the dataset increases. Dimension reduction is one alternative to ‘curb’ the combinatorial explosion of patterns by passing a reduced space to the analytical algorithms. The drawback, however, is the loss of patterns embedded

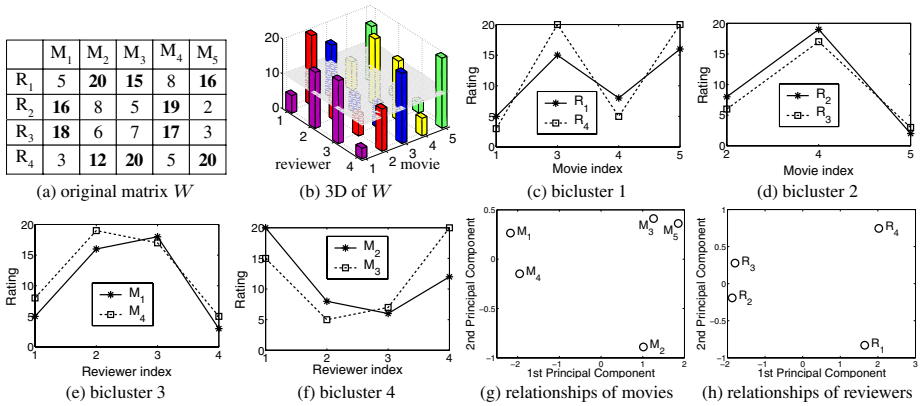


Fig. 1. The rating matrix for 4 reviewers and 5 movies, biclusters and PCA 2-dimensional space

in the subspace of the original space. This happens because most reduction techniques made use of distance or similarity measures over the full dimension, and therefore lack the mechanism to find the embedded patterns that are subtle but important.

In this paper, we introduce the concept of *data terrain* to visualize high-dimensional datasets while overcoming the limitations of subspace clustering and dimension reduction. Our proposal effectively reveals the relationship among subspace clusters, and allows the user to explore the data at different levels of details. We show, by means of real-world applications (e.g., biclusters, outliers, and frequent itemsets), how the data terrain can help discover generic patterns that can be utilized to effectively analyze the patterns embedded in the original space. Unfortunately, to find this data terrain under varying conditions proved to be NP-hard. Thus, our contribution in this paper includes the proposal of efficient techniques to find the data terrain. We next show a motivating example to illustrate the relevance of data terrains in analysis. We then introduce the resonance model in Section 3 and summarize our work in Section 4.

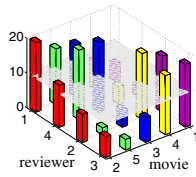
2 Motivating Example

We begin by introducing the concept of data terrain and show by means of an example, how it facilitates better data analysis; and why it is better than other techniques like biclustering and dimension reduction. Our example is based on the survey of popular movies. Fig. 1(a) shows the rating matrix W of 4 reviewers (R_i) on 5 movies (M_j), where each movie is rated on a scale of 1 to 20.

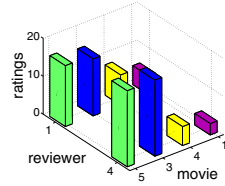
We first use biclustering to analyze the relationship between the reviewers and the movies. If requiring biclusters with at least 2 rows and columns, more than 10 biclusters can be discovered. Fig. 1(c) – (f) are the distinct biclusters found in this case. While these biclusters precisely characterized the reviewers’ ‘rating style’ on movies, there is too much redundancy in the solution for such a small dataset. In real-world situations where the dataset is much larger, it will take much longer before the an-

	M_2	M_5	M_3	M_4	M_1
R_1	20	16	15	8	5
R_4	12	20	20	5	3
R_2	8	2	5	19	16
R_3	6	3	7	17	18

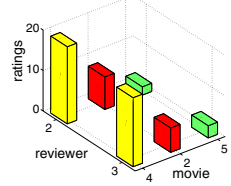
(a) matrix W' of reordered W



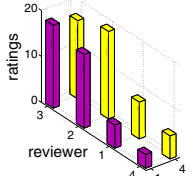
(b) 3D of W'



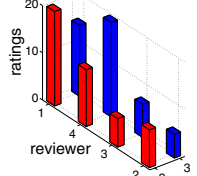
(c) bicluster 1



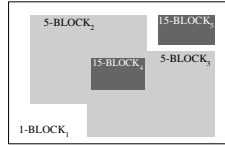
(d) bicluster 2



(e) bicluster 3



(f) bicluster 4



(g) BLOCKS and actual data



(h) Approximation of the actual terrain

Fig. 2. (a) – (f): Reordered rating matrix given by Fig. 1; (g) & (h): Illustrations of σ -BLOCKS

alyst is able to come to a conclusion. Worse, the analyst is likely to become confused on the real relationship between the reviewers and the movies if they were to work at this level of detail. The other approach would be to use dimension reduction techniques. We consider using a popular technique known as Principal Component Analysis (PCA). The PCA’s output is shown in Fig. 1(g) and (h). Again, the relationship between the reviewers and movies is revealed. However, as PCA runs across the full dimension of the data, we lose many subtle and local insights about reviewers and movies. For example, movies M_2 , M_4 and M_5 are perceived as different in Fig. 1(g). Yet, if we check back on our analysis using biclustering, we can see from Fig. 1(d) that they are actually quite similar if we consider the ratings from reviewer R_2 and R_3 . Thus, if only PCA is performed, we will not be able to arrive at this conclusion.

Interestingly, if we view W in 3D space, we can capture the relationships that both biclustering and PCA revealed. As Fig. 1(b) shows, a direct 3D ‘plot’ of W does not seem to reveal any interesting insights – but if we were to reorder W into W' as shown in Fig. 2(a) (and also Fig. 2(c) – (f); where every bicluster can be shown in this manner [1]), we have a 3D terrain of W' as depicted in Fig. 2(b). Notably, this terrain provides the insights that earlier requires both biclustering and PCA analysis.

To illustrate this, notice that any bicluster from Fig. 2(c) – (f) can be obtained by selecting some points from the ‘mountains’ and ‘plains’ in this terrain. At the same time, we can also make conclusions that would otherwise be obtained through PCA: (i) there are primarily two groups of movies and reviewers; (ii) M_3 and M_5 have higher similarity than M_2 despite being in the same group; and (iii) the ‘rating style’ of R_2 and R_3 is opposite that of R_1 and R_4 . Thus, the terrain captures both local and global relationships about the data in an intuitive and effective manner. Of course, real-world datasets are much more complex that result in more complicated terrains. Consequently, trying to discover such a terrain proved to be a NP-hard problem.

3 Discovering Data Terrains

Conceptually, moving from W to W' is simply the reordering of the matrix to form the ‘mountains’ and ‘plains’. Yet, this ordering on both dimensions can be difficult to achieve efficiently on massive datasets. To prove the hardness of this problem, we first give the following definitions. Let \mathcal{O} be a set of objects, where $o \in \mathcal{O}$ is defined by a set of attributes \mathcal{A} . Further, let w_{ij} be the magnitude of o_i over $a_j \in \mathcal{A}$. Then we can represent the relationship of all objects and their attributes in a matrix $W = (w_{ij})_{|\mathcal{O}| \times |\mathcal{A}|}$ for the weighted bipartite graph $G = (\mathcal{O}, \mathcal{A}, E, W)$, where E is the set of edges. Thus, discovering the ‘mountains’ transforms into the problem of evaluating subgraphs where the magnitude of all its edges are above some ‘altitude’, i.e., $w_{ij} \geq \sigma$. Formally, the concept of a ‘mountain’ in this data terrain is called a BLOCK.

Definition 1. Given a weighted bipartite graph G , a σ -BLOCK (or simply σ -B) is a subgraph $G' = (\mathcal{O}', \mathcal{A}', E', W')$ of G satisfying $w_{ij} \geq \sigma$ for any $i \in \mathcal{O}'$ and $j \in \mathcal{A}'$.

From Definition 1, σ -B can be intuitively viewed as a plane (or a transverse section) with a specified altitude σ that ‘cuts’ across W . In the case of Fig. 1(b), we set the plane at $\sigma=10$ to obtain two 10-Bs as shown in Fig. 2(b): $\{R_1, R_4, R_2\} \times \{M_2, M_5\}$ and $\{R_2, R_3\} \times \{M_4, M_1\}$. Therefore, a series of σ -Bs can be generated when considering planes with different σ values. Once this set of BLOCKs relevant to G is found, we can order them to find the data terrain.

Definition 2. Given a bipartite graph $G = (\mathcal{O}, \mathcal{A}, E, W)$ and a set of BLOCKs $\{B_1, B_2, \dots, B_k\}$ found from G , the terrain of W is two ordered sequences of \mathcal{O} and \mathcal{A} , such that these BLOCKs are placed consecutively in the reordered W .

It is interesting to note that sorting both dimensions, i.e., \mathcal{O} and \mathcal{A} , is an extension of sorting a single dimensional array to determine its distribution. However, sorting both dimensions simultaneously to get the 2-dimensional distribution is practically infeasible, i.e., finding the σ -BLOCKs by iteratively decreasing σ from the maximum value of W is NP-hard. In fact, finding a single σ -B is NP-hard.

Theorem 1. Finding the largest σ -BLOCK ($|\mathcal{O}'| \times |\mathcal{A}'|$) is NP-hard.

Proof. Our problem can be reduced from the *maximum edge biclique* [2], which is NP-complete. Details of this proof can be referred to [3].

Given the difficulty of finding σ -Bs, we seek alternative methods to discover the data terrain. Since our objective is to find the ‘mountains’ and ‘plains’ but *not* where they are on the terrain, then some approximation to the actual terrain (that is computationally efficient) should suffice. The insignificance of the specific locations of the ‘mountains’ and ‘plains’ can be demonstrated from Fig. 2(g) and (h), where the same set of insights are obtained from both figures. As this terrain is approximated, we called it the macro-view¹. To obtain the macro-view of a terrain for a dataset, we used a novel

¹ The complete work of this paper includes a *micro-view* of the data terrain. Together, they provide a complete solution for analysis of high-dimensional datasets. Due to space constraints, the reader is referred to [3] for the details.

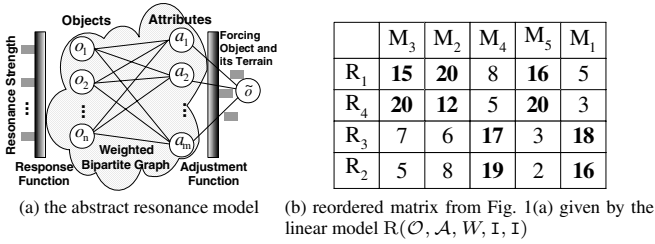


Fig. 3. The model and the effectiveness of the linear instance

model inspired by the physics of resonance. This resonance model is very efficient even on very large and high-dimensional datasets. Instead of checking every σ -B, we can simulate a resonance experiment by injecting a response function to elicit objects of interest to the analyst. Proofs of all theorems in this section are omitted and refer to [3].

3.1 The Model

To simulate a resonance phenomenon, we require a forcing object o , such that when an appropriate response function \mathbf{r} is applied, o will resonate to elicit those objects $\{o_i, \dots\} \subset \mathcal{O}$ in G , whose ‘natural frequency’ is similar to o . This ‘natural frequency’ represents the characteristics of both o and the objects $\{o_i, \dots\}$ who resonated with o when \mathbf{r} was applied. For the weighted bipartite graph $G = (\mathcal{O}, \mathcal{A}, E, W)$ and $W = (w_{ij})_{|\mathcal{O}| \times |\mathcal{A}|}$, this ‘natural frequency’ of $o_i \in \mathcal{O}$ is $\mathbf{o}_i = (w_{i1}, w_{i2}, \dots, w_{i|\mathcal{A}|})$. Since a one-dimensional array (or vector) can be sorted to obtain its own terrain, we also refer \mathbf{o}_i as the terrain of the object o_i . Likewise, the terrain of the forcing object o is defined as $\tilde{\mathbf{o}}_i = (w_1, w_2, \dots, w_{|\mathcal{A}|})$.

Put simply, if two objects of the same ‘natural frequency’ will resonate and therefore, should have a similar terrain. The evaluation of resonance strength between objects o_i and o_j is given by the response function $\mathbf{r}(\mathbf{o}_i, \mathbf{o}_j) : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$. We defined this function abstractly to support different measures of resonance strength. For example, one existing measure to compare two terrains is the well-known *rearrangement inequality theorem*, where $\mathbf{I}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n x_i y_i$ is maximized when the two positive sequences $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$ are ordered in the same way (i.e. $x_1 \geq x_2 \geq \dots \geq x_n$ and $y_1 \geq y_2 \geq \dots \geq y_n$) and is minimized when they are ordered in the opposite way (i.e. $x_1 \geq x_2 \geq \dots \geq x_n$ and $y_1 \leq y_2 \leq \dots \leq y_n$).

Notice if two vectors maximizing $\mathbf{I}(\mathbf{x}, \mathbf{y})$ are put together to form $M = [\mathbf{x}; \mathbf{y}]$ (in MATLAB format), we obtain the terrain. More importantly, all σ -Bs are immediately obtained from this terrain with the need to search every σ -B! This is why the model is efficient – it only needs to consider the resonance strength among objects once the appropriate response function is selected. For example, the response function \mathbf{I} is a suitable candidate to characterize the similarity of terrains of two objects. Likewise, $\mathbf{E}(\mathbf{x}, \mathbf{y}) = e^{-\left(\sum_{i=1}^n x_i y_i\right)}$ is also an effective response function.

To find the ‘mountains’ and ‘plains’, the forcing object o evaluates the resonance strength of every objects o_i against itself to locate a ‘best fit’ based on the contour of its terrain. By running this iteratively, those objects that resonated with o are discovered

and placed together to form the ‘mountains’ within the 2-dimensional matrix W . In the same fashion, the ‘plains’ are discovered by combining those objects that resonated weakly with o . This iterative learning process between o and G is outlined below.

Initialization. Set up o with a uniform distribution: $\tilde{\mathbf{o}} = (1, 1, \dots, 1)$; normalize it as $\tilde{\mathbf{o}} = \text{norm}(\tilde{\mathbf{o}})^2$; then let $k = 0$; and record this as $\tilde{\mathbf{o}}^{(0)} = \tilde{\mathbf{o}}$.

Apply Response Function. For each object $o_i \in \mathcal{O}$, compute the resonance strength $\mathbf{r}(\tilde{\mathbf{o}}, \mathbf{o}_i)$; store the results in a vector $\mathbf{r} = (\mathbf{r}(\tilde{\mathbf{o}}, \mathbf{o}_1), \mathbf{r}(\tilde{\mathbf{o}}, \mathbf{o}_2), \dots, \mathbf{r}(\tilde{\mathbf{o}}, \mathbf{o}_{|\mathcal{O}|}))$; and then normalize it, i.e., $\mathbf{r} = \text{norm}(\mathbf{r})$.

Adjust Forcing Object. Using \mathbf{r} from the previous step, adjust the terrain of o for all $o_i \in \mathcal{O}$. To do this, we define the adjustment function $\mathbf{c}(\mathbf{r}, \mathbf{a}_j) : \mathbb{R}^{|\mathcal{O}|} \times \mathbb{R}^{|\mathcal{O}|} \rightarrow \mathbb{R}$, where the weights of the j -th attribute is given in $\mathbf{a}_j = (w_{1j}, w_{2j}, \dots, w_{|\mathcal{O}|j})$. For each attribute a_j , $w_j = \mathbf{c}(\mathbf{r}, \mathbf{a}_j)$ integrates the weights from \mathbf{a}_j into o by evaluating the resonance strength recorded in \mathbf{r} . Again, \mathbf{c} is abstract, and can be materialized using the inner product $\mathbf{c}(\mathbf{r}, \mathbf{a}_j) = \mathbf{r} \bullet \mathbf{a}_j = \sum_i w_{ij} \cdot \mathbf{r}(\tilde{\mathbf{o}}, \mathbf{o}_i)$. Finally, we compute $\tilde{\mathbf{o}} = \text{norm}(\tilde{\mathbf{o}})$ and record it as $\tilde{\mathbf{o}}^{(k+1)} = \tilde{\mathbf{o}}$. We denote the resonance model as $\mathbf{R}(\mathcal{O}, \mathcal{A}, W, \mathbf{r}, \mathbf{c})$, where the instances of functions \mathbf{r} and \mathbf{c} can be either \mathbf{I} or \mathbf{E} .

Test Convergence. Compare $\tilde{\mathbf{o}}^{(k+1)}$ against $\tilde{\mathbf{o}}^{(k)}$. If the result converges, go to the next step; else apply \mathbf{r} on \mathcal{O} again (i.e., forcing resonance), and then adjust o .

Macro-View of Terrain. Sort the objects $o_i \in \mathcal{O}$ by the coordinates of \mathbf{r} in descending order; and sort the attributes $a_i \in \mathcal{A}$ by the coordinates of $\tilde{\mathbf{o}}$ in descending order.

3.2 Properties of the Model

The abstract view of the general model is given in Fig. 3(a). Depending on the response and adjustment function, the abstract model instantiates into different implementations. In practice, we have the linear model $\mathbf{R}(\mathcal{O}, \mathcal{A}, W, \mathbf{I}, \mathbf{I})$, and the non-linear model $\mathbf{R}(\mathcal{O}, \mathcal{A}, W, \mathbf{E}, \mathbf{E})$. We shall discuss some important properties of our model in this section. In particular, we show that the model gives a good approximation to the actual terrain, and that its iterative process converges quickly.

Approximation to Actual Terrain. Using the synthetic data from Fig. 2(g), we can see how well both implementations approximate the actual terrain. The linear and non-linear model converges to a precision of $\epsilon = 0.001$, i.e., once $\|\tilde{\mathbf{o}}^{k+1} - \tilde{\mathbf{o}}^k\| \leq \epsilon$, terminates. The reordered matrices are the same as Fig. 2(h). We then performed the same test on the movie-rating example. Result of linear model is shown in Fig. 3(b) and the non-linear model in Fig. 2(a). Obviously, $\mathbf{R}(\mathcal{O}, \mathcal{A}, W, \mathbf{E}, \mathbf{E})$ gives a better approximation, where the ‘mountains’ and ‘plains’ are easily distinguishable. Thus, we can conclude that the different instances of \mathbf{R} may give an approximate of the actual terrain. These conclusions are also empirically proven in [3].

Convergence. Since the resonance model is iterative, it is essential that it converges quickly to be efficient. Essentially, the model can be seen as a type of *discrete dynamical system* [4]. The convergence of linear and non-linear models is proven below.

Theorem 2. $\mathbf{R}(\mathcal{O}, \mathcal{A}, W, \mathbf{r}, \mathbf{c})$, where \mathbf{r}, \mathbf{c} are \mathbf{I} or \mathbf{E} , converges in limited iterations.

² $\text{norm}(\mathbf{x}) = \mathbf{x} / \|\mathbf{x}\|_2$, where $\|\mathbf{x}\|_2 = (\sum_{i=1}^n x_i^2)^{1/2}$ is 2-norm of vector $\mathbf{x} = (x_1, \dots, x_n)$.

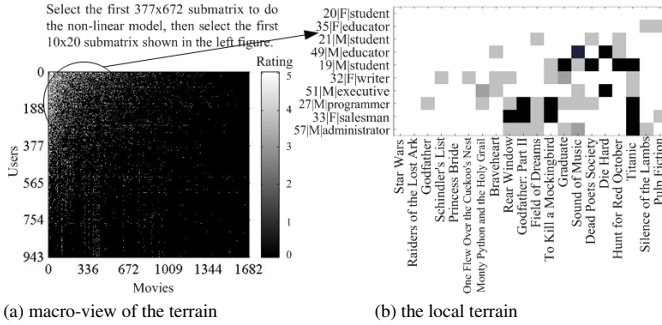


Fig. 4. A case of data analysis by macro-view of the terrain: visualization of MovieLens data (with 943 users, 1,682 movies and 100,000 ratings on the scale of 1 to 5. It is available at <http://www.grouplens.org>), and analysis of its local terrain. We obtain the macro-view using $R(\mathcal{O}, \mathcal{A}, W, I, I)$ to get the terrain in (a). It is possible that the user is not satisfied with an overview of the dataset and might be interested in further analysis. A case of ‘zoom in’ on the ‘crowded’ but a local terrain of the macro-view is shown in (b).

In practice, the model is very efficient because we are only interested in the convergence of orders of coordinates in $\tilde{\sigma}^k$ and \mathbf{r}^k . With k iterations, the complexity is $\mathcal{O}(k \times |\mathcal{O}| \times |\mathcal{A}|)$. In our experiments, our model converges within 50 iterations even on the non-linear configurations giving a time complexity of $\mathcal{O}(|\mathcal{O}| \times |\mathcal{A}|)$. In all cases, the complexity is sufficiently low to efficiency handle large datasets.

Average Inter-resonance Strength $\frac{1}{\binom{k}{2}} \sum_{\substack{i \neq j \in \mathcal{O}' \\ |\mathcal{O}'|=k}} \mathbf{r}(\mathbf{o}_i, \mathbf{o}_j)$ **among Objects.** Theorem 3 is in fact an optimization process to find the best k objects, whose average inter-resonance strength is the largest among any subset of k objects. Next, we exploit this property to unveil the relationship between the macro-view of the data terrain and the biclusters.

Theorem 3. *Given the macro-view terrain W' , the average inter-resonance strength $\frac{1}{\binom{k}{2}} \sum_{1 \leq i \neq j \leq k} \mathbf{r}(\mathbf{o}_i, \mathbf{o}_j)$ of the first k objects, w.r.t. the resonance strength with \mathbf{o} , is largest for any subset with k objects.*

Approximation to Maximum Edge Biclique (MEB). The non-linear configuration of our model, i.e., $R(\mathcal{O}, \mathcal{A}, W, E, E)$ has such capability. Details refer to [3].

3.3 Real World Examples

A demonstration of how a macro-view of the data terrain can help the user in analysis is shown by a real-world case in Fig. 4. Next we show how it can have applications in data mining for finding frequent itemsets and biclustering in theory. All empirical evidences refer to [3].

Finding Frequent Itemsets. A transaction dataset can be constructed as a matrix, where each transaction is an object, and each item is an attribute whose value w_{ij} in $W_{|\mathcal{O}| \times |\mathcal{A}|}$ is 1 if the j -th item occurs in the i -th record, and 0 otherwise. We therefore have the following that relates frequent itemsets and BLOCKs.

Theorem 4 (Frequent Itemsets and BLOCKS). *A frequent itemset is the attribute set of a 1-BLOCK, and its support is the number of objects in the BLOCK.*

Discovering Biclusters. A popular measure for biclusters [5] is defined as Eqn. (1). The residue $H(W)$ of given a matrix W is a δ -bicluster if $H(W) \leq \delta$.

$$H(W) = \frac{1}{mn} \sum_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} (w_{ij} - w_{iJ} - w_{Ij} + |W|)^2 \tag{1}$$

where $w_{iJ} = \frac{1}{n} \sum_{j=1}^n w_{ij}$, $w_{Ij} = \frac{1}{m} \sum_{i=1}^m w_{ij}$, and $|W| = \frac{1}{mn} \sum_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} w_{ij}$.

Theorem 5 (Bicluster and Average Resonance Strength of Macro-View Terrains). *Given a matrix $W = (w_{ij})_{m \times n}$, where \mathcal{O} are the rows and \mathcal{A} are columns, we have the inverse relation of the average inter-resonance strength and $H(W)$ as follows*

$$H(W) = \|W\|^2 + |W|^2 - \frac{1}{n} \bar{r}(W) - \frac{1}{m} \bar{r}(W^T) \tag{2}$$

where $\|W\| = \sqrt{\frac{1}{mn} \sum_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} w_{ij}^2}$, $\bar{r}(W) = \frac{1}{\binom{m}{2}} \sum_{\substack{1 \leq i, j \leq m \\ i \neq j}} I(\mathbf{w}_{i\cdot}, \mathbf{w}_{j\cdot})$ is the average inter-resonance strength among $\mathbf{w}_{i\cdot}$, and $\bar{r}(W^T) = \frac{1}{\binom{n}{2}} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} I(\mathbf{w}_{\cdot i}, \mathbf{w}_{\cdot j})$ is the average inter-resonance strength among $\mathbf{w}_{\cdot j}$, and $\mathbf{w}_{i\cdot}$ is the i -row vector of W with $\mathbf{w}_{\cdot j}$ the j -column vector of W .

It can be interpreted as follows. Since $\|W\|$ and $|W|$ are sum of W in different forms, we can consider them as fixed constant. If the average inter-resonance strength of W and W^T , i.e., $\bar{r}(W)$ and $\bar{r}(W^T)$, is higher, then $H(W)$ is lower and thus, W behaves like a bicluster. For $R(\mathcal{O}, \mathcal{A}, W, I, I)$ and $R(\mathcal{A}, \mathcal{O}, W^T, I, I)$, we conclude that if we select the first k rows and columns of W with large resonance strength $r(\mathbf{o}_i, \tilde{\mathbf{o}})$ to form W' , it is straightforward that we will have a smaller $H(W')$ and thus, W a bicluster.

4 Summary

In this paper, we proposed the data terrain as a means to visualize and analyze high-dimensional datasets. With this terrain, patterns in subspaces can be visualized and analyzed. We provided a novel solution to obtain the the macro-view of a terrain efficiently, and demonstrated its real-world application.

References

1. Liu, J., Wang, W.: Op-cluster: Clustering by tendency in high dimensional space. In: Proceedings of ICDM, Melbourne, Florida (2003) 187
2. Peeters, R.: The maximum edge biclique problem is NP-complete. *Disc. App. Math.* **131** (2003) 651–654
3. Li, W., Ong, K.L., Ng, W.K.: Visual terrain analysis of high dimensional datasets. Technical Report (www.deakin.edu.au/leong/tr0406) (TRC04/06), Deakin University (2005)
4. Sandefur, J.T.: *Discrete Dynamical Systems*. Oxford: Clarendon Press (1990)
5. Cheng, Y., Church, G.M.: Biclustering of expression data. In: Proceedings of the 8th International Conference on Intelligent System for Molecular Biology. (2000)

An Auto-stopped Hierarchical Clustering Algorithm for Analyzing 3D Model Database

Tian-yang Lv^{1,2}, Yu-hui Xing², Shao-bing Huang¹, Zheng-xuan Wang²,
and Wan-li Zuo²

¹ College of Computer Science and Technology, Harbin Engineering University,
Harbin, China

² College of Computer Science and Technology, Jilin University, Changchun, China
raynor1979@163.com, wanli@mail.jlu.edu.cn

Abstract. In the research of shape-based 3D model retrieval, the analysis and classification of 3D model database is an important topic for improving the retrieval performance. However, it encounters difficulties due to lack of valuable prior knowledge and the semantic gaps exist in 3D model retrieval. The paper proposes a new auto-stopped hierarchical clustering algorithm overcome these problems, which combines outlier detection with clustering. The Princeton Shape Benchmark along with 2 data sets from UCI is employed to evaluate the performance of the algorithm. And the new algorithm outperforms other auto-stopped algorithms and obtains better classification of 3D model database.

Keywords: shape-based 3D model retrieval; clustering; outlier detection.

1 Introduction

With the proliferation of 3D models and their wide spread through internet, 3D model retrieval, especially shape-based 3D model retrieval, becomes a new emerging research field ^[1]. However, as an important subtopic, the analysis and organization of the 3D model databases encounters difficulties due to lack of the valuable domain knowledge. For instance, little is known about the number of models' classes. Moreover, the two-level semantic gaps exist in 3D model retrieval: one is the gap between the shape of model and its feature, which means models with similar shape have great different feature; the other is the gap between the shape of model and its meaning in real-life, which results in the mistakes in manually classifying the 3D model database.

The paper explores the application of the clustering techniques in analyzing 3D model database. And the clustering result is treated as the classification of 3D model database, since models of the same cluster have similar feature.

The topic has not been thought much in the previous works. For example, it is very difficult to pre-decide an appropriate number of final clusters k for 3D model database, while k is required by many traditional clustering algorithms, such as the hierarchical clustering algorithms CURE ^[2] and the partitioning algorithm K-means.

Thus, the paper proposes an auto-stopped hierarchical clustering algorithm, which integrates a new outlier mining method in clustering and cancels the parameter k . It is based on the following observation: the distances among data or clusters not only show their similarity degree, but also demonstrate the dissimilarity. With the pro-

gressing of clustering, the dissimilarity $D(C_{NN-A}, C_{NN-B})$ between the two most similar clusters C_{NN-A} and C_{NN-B} at present is increasing. And the clustering should stop at the moment when C_{NN-A} and C_{NN-B} are so diverse from each other. The outlier-mining process can provide that suitable “diverse degree” since outliers are detected according to their “great difference” from the others.

The rest of paper is organized as follows: after introducing related works in section 2, section 3 proposes the new clustering algorithm; section 4 gives the experimental result; finally, section 5 summarizes the paper.

Table 1. Important Notations

Notation	Description
N	Total number of Data
M	Dimensionality of Data
C_i	The i th Cluster
$D(C_i, C_j)$	Distance between C_i and C_j

2 Related Works

CURE algorithm ^[2] employs the novel concept of *representative* to represent a cluster and r *representatives* are shrunk towards the cluster’s centroid by a fraction α before computing clusters’ distance to avoid noise. However, CURE needs the parameter k and does not consider clusters’ density in merging decisions.

Some researches try to make clustering algorithm optimally estimate k . [3] proposes a method based on dissimilarity increment. But it is short at handling outlier and detecting clusters with complex shape, like the linearly inseparable datasets.

To achieve the property of rotation invariance, [1] states a method using spherical harmonic transformation on voxel descriptors of 3D model. Its overview is: first, the 3D model is projected into a $2R \times 2R \times 2R$ voxel grid and set the corresponding value of a voxel 1, if it contains point of polygonal surface, otherwise set the value of 0; then, normalize the model with translation and scale; thus, for each sphere with the radius r , the spherical function of a 3D model can be defined as:

$$f_r(\theta, \varphi) = \text{Voxel}(r \sin(\theta) \cos(\varphi) + R, r \cos(\theta) + R, r \sin(\theta) \sin(\varphi) + R) \quad (1)$$

where $\theta \in [0, \pi], \varphi \in [0, 2\pi]$ and $r \in [0, R]$. And for each spherical harmonic function f_r can be decomposed as the sum of different frequencies, like:

$$f_r(\theta, \varphi) = \sum_{l=0}^{B-1} f_r^l(\theta, \varphi), f_r^l(\theta, \varphi) = \sum_{m=-l}^l a_{l,m} Y_l^m(\theta, \varphi) \quad (2)$$

Where $Y_l^m(\theta, \varphi)$ is the harmonic homogeneous polynomial of l . Combining the signature $\{\|f_r^0\|, \|f_r^1\|, \dots\}$ for f_r with different r , the shape descriptor for the 3D model is obtained, whose dimensionality depends on B and R with R usually equals 32.

3 An Auto-stopped Hierarchical Clustering Algorithm

Based on the traditional hierarchical clustering process, the Auto-Stopped Clustering Algorithm using Representatives ASCAR shows its uniqueness in three aspects: (1) adopts a new distance-based outlier detection method before clustering to detect outliers and exclude their disturbance for the clustering process; (2) employs the representatives and considers clusters' density in computing clusters' distance; (3) stops clustering automatically according to the dissimilarity reflected by the outliers.

3.1 The Outlier Detection Method Based on Even-Distribution Pattern

The basic idea of distance-based outlier detection method is: if the distances between data a and most other data are larger than the threshold D_{out} , a is an outlier [4]. It is critical but usually difficult to decide an appropriate D_{out} . This method also ignores the local distribution feature of one data.

The new method decides D_{out} according to the even distribution pattern of data. It is a very useful reference, since clusters and outliers exist only if the real-life data distribute unevenly. In that case, the distances \bar{D}_{NN} between each data and its nearest neighbor are the same. \bar{D}_{NN} is approximately decided according to equation 3, where $a_{max}^{(i)}$ and $a_{min}^{(i)}$ is the maximum and the minimum of all data's i th-dimension. And $D_{out} = \bar{D}_{NN} / \beta$, where β is a parameter to describe the diversity of the realistic distribution situation from the even pattern.

$$\bar{D}_{NN} = \sqrt{\sum_{i=1}^M ((a_{max}^{(i)} - a_{min}^{(i)}) / \sqrt[M]{N})^2} \tag{3}$$

Factor ξ is adopted to evaluate the local distribution feature of a data. For data a , $\xi(a) = D_{NN}(a) / D_{NN}(b)$, where $D_{NN}(a)$ is the distance between a and its nearest-neighbor and so is $D_{NN}(b)$. The value of $\xi(a)$ shows the isolation degree of a from its neighbors. Special method is used for very similar or duplicate data. And, the equation of $\xi(a)$ is:

$$\xi(a) = \begin{cases} D_{NN}(a) / D_{NN}(b) & \text{if } (D_{NN}(b) > 10^{-4}) \\ 1 & \text{else} \end{cases} \tag{4}$$

Therefore, the outlier evaluation criterion is stated as follows:

Data a is an outlier, if $D_{NN}(a) * \xi(a) > (\frac{\sqrt{\sum_{i=1}^M ((a_{max}^{(i)} - a_{min}^{(i)}) / \sqrt[M]{N})^2}}{\beta})$

Since outliers are extremely far away from the others while the normal are relatively near to each other, a method is proposed to decide the appropriate β :

(1) name β_{Step} as the *step length* and $\beta_{Step} = D_{NN}(a_{forest}) * \xi(a_{forest}) / \bar{D}_{NN}$, where a_{forest} satisfies $D_{NN}(a_{forest}) * \xi(a_{forest}) \geq D_{NN}(b) * \xi(b)$ for any b ; (2) observe the increasing speed V of the detected outlier number n_{out} under different value of β ,

viz. $V = \nabla n_{out} / \nabla \beta = \nabla n_{out} / \beta_{Step}$, where $\beta = l \times \beta_{Step}$ and $l = \{1, 2 \dots\}$, call l the *step Num.*; (3) if V reaches its first peak when $l_i \times \beta_{Step}$, $\beta = (l_i - 1) \times \beta_{Step}$.

3.2 The Computation of the Clusters' Distance

The algorithm ASCAR adopts *representative* from CURE algorithm to improve clustering performance. But ASCAR excludes the influence of “noise” by adopting a professional outlier mining method without using the parameter α . And ASCAR also considers the cluster’s density in deciding whether clusters should be merged.

The algorithm decides the distance $D(C_i, C_j)$ between C_i and C_j according to two factors: first, the distance $D_{min}(C_i, C_j)$ of the nearest *representatives* coming from C_i and C_j respectively; second, the factor δ measuring the change of cluster’s density *Den.* The density of C_i or C_j approximately equals the average distances among its *representatives*. For the new-borne cluster C_{new} created by merging C_i and C_j , $Den(C_{new}) = D_{min}(C_i, C_j)$. Then, $\delta(C_i)$ is defined as follows and so is $\delta(C_j)$:

$$\delta(C_i) = \begin{cases} Den(C_i) / Den(C_{New}) & \text{if } (Den(C_i) > Den(C_{New})) \\ Den(C_{New}) / Den(C_i) & \text{otherwise} \end{cases} \quad (5)$$

Since it is impossible to compute the density of the cluster with only one data, define $D(C_i, C_j) = D_{min}(C_i, C_j)$ in that case. And the way to compute $D(C_i, C_j)$ is:

$$D(C_i, C_j) = \begin{cases} D_{Min}(C_i, C_j) \times (\delta(C_i) + \delta(C_j)) / 2 & \text{if } (n_i > 1) \ \& \ (n_j > 1) \\ D_{Min}(C_i, C_j) & \text{otherwise} \end{cases} \quad (6)$$

Factor δ reflects the influence of cluster’s density on merging decision. That is, the bigger the difference between the density of C_i or C_j with that of C_{new} , the less possibility for C_i and C_j to be merged.

3.3 Automatic Stop

Without user-specified condition to stop clustering, it is necessary to extract this information from the processed data. As stated in former parts, it is a suitable opportunity to stop clustering if the clusters to be merged are too dissimilar. We propose D_{out} as the dissimilarity threshold to decide this opportunity for two reasons: (1) D_{out} is used to detect outliers, while the major characteristic of outliers is their dissimilarity from the others; (2) D_{out} is decided according to the even distribution pattern, which is also a useful reference for cluster analysis since the existence of clusters shows the diversity of the realistic distribution situation from the even pattern.

And the stop criterion of clustering is :

Suppose C_{NN-A} and C_{NN-B} are the most similar clusters at present, stop clustering if $D(C_{NN-A}, C_{NN-B}) > D_{out}$.

3.4 Complexity Analysis and Overview of ASCAR

The complexity of traditional hierarchical algorithm is $O(N^2)$. Since ASCAR is constructed on the traditional method, it is only necessary to analyze the complexity of

each change. The complexity increases by $O(N)$ to perform one more scan to detect outliers. The complexity increases by $O((r*n_i+r^2)*(N-k))$ at most in computing clusters' distance. Thus, the complexity increases by $O((r*n_i+r^2)*(N-k) +N)$ in total. Since $r^2 < N$ in most cases, the complexity of ASCAR equals $O(N^2)$.

The overview of the proposed algorithm ASCAR is listed in Figure 1.

```

Algorithm ASCAR (  $r, \beta$  )
1. { Read all data and decide vector  $a_{max}$  and  $a_{min}$ ;
2.  Treat each data as a separate cluster;
3.  Compute each cluster's nearest-neighbor;
4.  Determine the value of  $\beta_{Step}$ ;
5.   $D_{out} = \mathbf{outlier}(a_{max}, a_{min}, \beta_{Step})$ ;
6.  Name the nearest clusters at present as  $C_{NN-A}, C_{NN-B}$ ;
7.  while ( $D(C_{NN-A}, C_{NN-B}) \leq D_{out}$ )
8.  { Merge clusters  $C_{NN-A}$  and  $C_{NN-B}$ ;
9.    Update  $C_{NN-A}$  and  $C_{NN-B}$ ; }
10. } //End of ASCAR
    
```

Fig. 1. The Auto-stopped Hierarchical Clustering Algorithm ASCAR

4 Experiment and Analysis

The evaluation of the new algorithm is undertook in two aspects: first, the real-life datasets *Iris* and *Wine* of UCI Machine Learning Repository [5] are adopted; then ASCAR is applied in analyzing the Princeton Shape Benchmark [8].

4.1 Data Sets of UCI

The criterions *Entropy* and *Purity* of [6] are adopted to measure the clustering results' quality for *Iris* and *Wine* datasets. And the better the clustering result, the smaller is

Table 2. Overview of the clustering results of ASCAR and other algorithms

	Dataset	Parameters	k	Entropy	Purity	n_{out}
ASCAR	Iris	$\beta=2.2, r=5$	5	0.3542	0.8121	5
	Wine	$\beta=1.8, r=4$	14	0.1837	0.8864	3
Frozen	Iris	$A=4.0$	2	0.4206	0.6667	--
	Wine	$A=0.5$	13	0.4998	0.7247	--
DBScan	Iris	$\epsilon=0.7, MPts=3$	2	0.4077	0.6867	--
	Wine	$\epsilon=35, MPts=3$	6	0.5866	0.6798	--

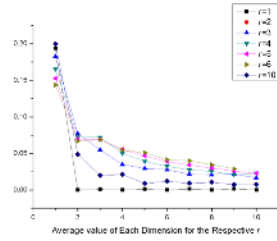
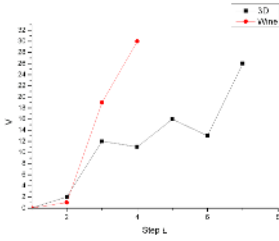


Fig. 2. Changes of V with the Increase of l **Fig. 3.** Average of each dimension f

the *Entropy* and the bigger is the *Purity*. The clustering performance of ASCAR is listed in Table 2 along with the detected number of outlier. To be more persuasive, Table 2 also gives the best clustering results of DBScan^[7] and Frozen. Figure 2 shows the change of V with the increasing of *step num.* l .

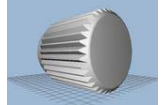
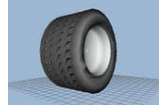

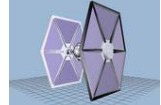
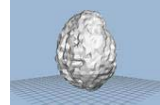
4.2 Princeton Shape Benchmark

The feature extraction method of section 2 with $R=32$ and $B=10$ is applied to obtain the shape feature from 3D models. Obviously, the model feature with the dimensionality 320 will greatly reduce the clustering performance. Figure 3 shows that the first element of the transformation result for each $f_r(\theta, \varphi)$ plays the most important role in distinguishing model. Therefore, we select that element and obtain the shape feature with $M=32$. In experiment, the Euclidean distance is adopted.

Table 3. Cluster’s detail of C_{90} and C_{160}

No	Cluster’s Detail				
C_{90}					
C_{160}					

Table 4. Part of the detected outliers and the respective value of *step num l*

M741 (L=1)	M737(L=2)	M416(L=2)	M1401(L=3)	M286(L=3)
				

Since there is no valuable knowledge of the classification of the models in PSB, we have to list the details of the result cluster. When $r=4$ and $\beta=0.8*5$, ASCAR obtains 160 clusters with the smallest size of 2. Due to space limit, Table 3 just lists the details of C_{90} and C_{160} . Comparing to the manual classification result of PSB, ASCAR achieves very similar classes' number. But, ASCAR clusters the models with similar shape together no matter what real-life meaning they represent, especially if the feature extraction method satisfies the request that models with similar shape have similar feature. However, this cannot be not always satisfied and clustering mistakes can be observed in Table 3. Table 4 lists part of the detected outliers along with *step num l*, under which they are pruned.

We also applied the auto-stopped algorithms DBScan and Frozen in analyzing PSB. Under all possible value of parameters, Frozen algorithm obtains over 1200 clusters, while DBScan tends to obtain a little huge clusters. For instance, when $\epsilon=0.2$ and $MPts=2$, DBScan gets 66 clusters with $n_0=715$, $n_1=1028$, $n_2, n_3...n_7=2$, and $n_8, n_9...n_{65}=1$. Obviously, these results are not acceptable as a classification of the database.

5 Conclusion

To analyze the 3D model database, the paper proposes a new strategy that integrates outlier detection with clustering and introduces an auto-stopped hierarchical clustering algorithm ASCAR. Experimental results show ASCAR's good performance in clustering the Princeton Shape Benchmark and 2 datasets from UCI. The future works will concentrate on the study of using the representations of the clustering result to establish the index of 3D model database.

Acknowledgements

This work is sponsored by the Natural Science Foundation of China under grant number 60373099 and the Natural Science Research Foundation of Harbin Engineering University under the grant number HEUFT05007.

References

1. T.Funkhouser, et al. A Search Engine for 3D Models. ACM Transactions on Graphics.22 (1), (2003) 85-105.
2. S. Guha, R. Rastogi, K. Shim: CURE: an Efficient Clustering Algorithm for Large Database. In: Laura M. Haas and Ashutosh Tiwary, eds. Proceedings of the ACM SIGMOD Conference on Management of Data. Seattle, Washington: ACM Press (1998) 73-84.

3. Ana L.N. Fred, José M.N. Leitão: A new Cluster Isolation criterion Based on Dissimilarity Increments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, VOL. 25, No. 8 (August 2003) 944-958
4. Edwin M. Knorr, Raymond T. Ng: Finding Intensional Knowledge of Distance-Based outliers. In: *Proceedings of the 25th Very Large Data Bases conference*. Edinburgh, Scotland (1999) 211 - 222
5. Hettich, S. & Blake, C.L. & Merz, C.J. (1998). UCI Repository of machine learning databases [<http://www.ics.uci.edu/~mlearn/MLRepository.html>]. Irvine, CA: University of California, Department of Information and Computer Science.
6. Ying Zhao, George Karypis: Criterion Functions for Document Clustering: Experiment and Analysis. Technical Report #01-40, University of Minnesota (2001) 1 – 40
7. M. Ester, H.-P. Kriegel, J. Sander, and X. Xu: A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In: *Proceedings of 2nd International Conference on Knowledge Discovery and Data Mining*, Portland, OR (1996) 226-231
8. Philip Shilane, Patrick Min, Michael Kazhdan, and Thomas Funkhouser: The Princeton Shape Benchmark. *Shape Modeling International*, Genova, Italy, (June 2004)

A Comparison Between Block CEM and Two-Way CEM Algorithms to Cluster a Contingency Table

Mohamed Nadif¹ and Gérard Govaert²

¹ LITA EA3097, Université de Metz, Ile du Saulcy, 57045 Metz, France
mohamed.nadif@iut.univ-metz.fr

² HEUDIASYC, UMR CNRS 6599, Université de Technologie de Compiègne,
BP 20529, 60205 Compiègne Cedex, France
gerard.govaert@utc.fr

Abstract. When the data consists of a set of objects described by a set of variables, we have recently proposed a new mixture model which takes into account the block clustering problem on the both sets and have developed the *block CEM* algorithm. In this paper, we embed the block clustering problem of contingency table in the mixture approach. In using a Poisson model and adopting the classification maximum likelihood principle we perform an adapted version of block CEM. We evaluate its performance and compare it to a simple use of CEM applied on the both sets separately. We present detailed experimental results on simulated data and we show the interest of this new algorithm.

1 Introduction

Cluster analysis is an important tool in a variety of scientific areas such as pattern recognition, information retrieval, micro-array, data mining, and so forth. Although many clustering procedures such as hierarchical clustering, k -means or self-organizing maps, aim to construct an optimal partition of objects or, sometimes, of variables, there are other methods, called block clustering methods, which consider simultaneously the two sets and organize the data into homogeneous blocks.

A wide variety of procedures have been proposed for finding patterns in data matrices. These procedures differ in the pattern they seek, the types of data to which they apply, and the assumption on which they rely. Let us mention the works of Hartigan (1975), Bock (1979), Garcia and Proth (1986), Marchotorchino (1987), Govaert (1983, 1995), Arabie and Hubert (1990), Duffy and Quiroz (1991) and Ritschard et al. (2001) who have proposed some algorithms dedicated to different kinds of matrices.

These last years, block clustering (also called biclustering) has become an important challenge in data mining context. In the text mining field, Dhillon (2001) has proposed a spectral block clustering method by exploiting the duality between rows (documents) and columns (words). In the analysis of micro-array

data where data are often presented as matrices of expression levels of genes under different conditions, block clustering of genes and conditions has permitted to overcome the problem of the choice of similarity on the both sets found in conventional clustering methods (Cheng and Church, 2000). Also, these kinds of methods have practical importance in a wide of variety of applications such as text and market basket data analysis. Typically, the data that arises in these applications is arranged as a two-way contingency or co-occurrence table.

In this paper, we will focus on these kinds of data. The data which we consider is noted \mathbf{x} ; it is a $r \times s$ data matrix defined by $\mathbf{x} = \{(x_{ij}); i \in I, j \in J\}$, where I is a categorical variable with r categories and J a categorical variable with s categories. In exploiting the duality between I and J , we will study the block clustering problem in embedding it in the mixture approach. We will propose a *block mixture model* which takes into account the block clustering situation and perform an innovative co-clustering algorithm. This one is based on the alternated application of Classification EM (Celeux and Govaert, 1992) on intermediate data matrices. To propose this algorithm, we set this problem in the classification maximum likelihood (CML) approach (Symons, 1981). This paper deals to compare block CEM and two-way CEM, i.e. CEM applied separately on I and J . Results on simulated data are given, confirming that block CEM gives much better performance than two-way CEM.

The paper is organized as follows. In Section 2, we give the necessary background CML approach and we describe the CEM algorithm and its steps when the data is arranged as a two-way contingency. In Section 3, we start by recalling our block mixture model and we describe the block CEM algorithm. In order to compare two-way CEM and block CEM, in Section 4, we perform numerical Monte Carlo simulations. A final section summarizes and indicates the recommended algorithm.

2 Mixture Model and Clustering

For convenience, we represent a partition of I into g clusters by $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_g)$ where \mathbf{z}_i , which indicates the component of the row i , is represented by $\mathbf{z}_i = (z_{i1}, \dots, z_{ig})$ with $z_{ik} = 1$ if row i is in cluster k and 0 otherwise. Then, the k th cluster corresponds to the set of rows i such that $z_{ik} = 1$. We will use similar notation for a partition \mathbf{w} into m clusters of the set J . In the following, to simplify the notation, the sums and the products relating to rows, columns or clusters will be subscripted respectively by letters i, j or k without indicating the limits of variation, which will be thus implicit. Thus, for example, the sum \sum_i stands for $\sum_{i=1}^r$ or $\sum_{i,j,k,\ell}$ stands for $\sum_{i=1}^r \sum_{j=1}^s \sum_{k=1}^g \sum_{\ell=1}^m$.

2.1 CML Approach and the CEM Algorithm

In the model-based clustering (see for instance (McLachlan and Peel, 2000)), it is assumed that the data are generated by a mixture of underlying probability

distributions, where each component k of the mixture represents a cluster. Thus, the density of the observed data \mathbf{x} is expressed as

$$f(\mathbf{x}; \boldsymbol{\theta}) = \prod_i \sum_k \pi_k \varphi_k(\mathbf{x}_i; \alpha_k) \tag{1}$$

where $\boldsymbol{\theta} = (\pi_1, \dots, \pi_g, \alpha_1, \dots, \alpha_g)$, (π_1, \dots, π_g) are the mixing proportions and $(\alpha_1, \dots, \alpha_g)$ are the parameters of the density components φ_k .

The clustering problem can be studied under mixture model using two different approaches: the maximum likelihood (ML) approach and the classification maximum likelihood (CML) approach (Symons, 1981). In this paper we focus on the second approach.

The ML approach estimates the parameters of the mixture and the partition is derived from these parameters using the maximum a posteriori principle (MAP). In the CML, the partition is added to the parameters to be estimated. The CML approach consists in estimating the parameters of the mixture and the partition. The maximum likelihood estimation of these new parameters leads to optimize in $\boldsymbol{\theta}$ and \mathbf{z} the complete data log-likelihood

$$L_C(\mathbf{z}, \boldsymbol{\theta}) = L(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z}) = \sum_{i,k} z_{ik} \cdot g(p_k \varphi_k(\mathbf{x}_i; \alpha_k)).$$

This optimization can be done by the Classification EM (CEM) algorithm (Celex and Govaert, 1992), a variant of EM (Dempster, Laird and Rubin, 1977), which converts the posterior probabilities t_{ik} 's to a discrete classification in a C-step before performing the M-step.

2.2 Application to Contingency Table

In this situation, the contingency table \mathbf{x} is a $r \times s$ data matrix defined by $\mathbf{x} = \{(x_{ij}); i \in I, j \in J\}$, where I and J are categorical variables with r and s categories. The sum of each row i will be denoted $x_{i\cdot}$. Thus, if we note $\boldsymbol{\theta} = (\pi_1, \dots, \pi_g, \alpha_{11}, \dots, \alpha_{gs})$ the parameter of the model and φ is the multinomial distribution of the k -th component, the log-likelihood (up to a constant) can be written as $L(\boldsymbol{\theta}; \mathbf{x}) = \sum_i \cdot g \sum_k \pi_k \alpha_{k1}^{x_{i1}} \dots \alpha_{ks}^{x_{is}}$, and the complete data log-likelihood as $L(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z}) = \sum_{i,k} z_{ik} \left(\cdot \pi_k + \sum_j x_{ij} \cdot g \alpha_{kj} \right)$.

In clustering context, the use of the mixture model deals to find the component from which each row arises. The CEM algorithm allows us to achieve this goal and the different steps of CEM in this situation are

- E-step: compute the posterior probabilities $t_{ik}^{(c)} \propto \pi_k \alpha_{k1}^{x_{i1}} \dots \alpha_{ks}^{x_{is}}$;
- C-step: the k th cluster of $\mathbf{z}^{(c+1)}$ is defined with $z_{ik}^{(c+1)} = 1$ if $k = \operatorname{argmax}_{k=1, \dots, g} t_{ik}^{(c)}$ and $z_{ik}^{(c+1)} = 0$ otherwise;
- M-step: by standard calculations, one arrives at the following re-estimations parameters $\pi_k^{(c+1)} = \frac{n_k^{(c+1)}}{r}$ and $\alpha_{kj}^{(c+1)} = \frac{x_{kj}}{x_k}$ where $n_k^{(c+1)}$ is the cardinality of the k th cluster of $\mathbf{z}^{(c+1)}$, $x_{kj} = \sum_i z_{ik}^{(c+1)} x_{ij}$ and $x_k = \sum_j x_{kj}$.

Having found the estimate of the parameters and noting $f_{kj} = \frac{x_{kj}}{x_{..}}$ where $x_{..} = \sum_{i,j} x_{ij}$, we can show that, when the proportions are fixed, the maximization of $L(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z})$ is equivalent to the maximization of the mutual information $I(\mathbf{z}, J) = \sum_{k,j} f_{kj} \cdot g \frac{f_{kj}}{f_{k.} f_{.j}}$ and approximately equivalent to the maximization of the chi-square criterion $\chi^2(\mathbf{z}, J) = x_{..} \sum_{k,j} \frac{(f_{kj} - f_{k.} f_{.j})^2}{f_{k.} f_{.j}}$. Hence the use of the both criteria $\chi^2(\mathbf{z}, J)$ and $I(\mathbf{z}, J)$ supposes implicitly that the data arise from a mixture of multinomial distributions. To tackle the block clustering problem, we can obviously use the CEM on I and J separately (noted 2CEM) but unfortunately it is unaware of the correspondence between I and J . It will be seen later that this process is ineffective to detect homogeneous blocs.

3 Block Mixture Model for Contingency Table

To study the block clustering problem, we have extended (Govaert and Nadif, 2003) the mixture model to propose a block mixture model defined by the following probability density function

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{(\mathbf{z}, \mathbf{w}) \in \mathcal{Z} \times \mathcal{W}} \prod_i \pi_{z_i} \prod_j \rho_{w_j} \prod_{i,j} \varphi(x_{ij}; \boldsymbol{\alpha}_{z_i w_j})$$

where $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\rho}, \boldsymbol{\alpha}_{11}, \dots, \boldsymbol{\alpha}_{gm})$, $\boldsymbol{\pi} = (\pi_1, \dots, \pi_g)$ and $\boldsymbol{\rho} = (\rho_1, \dots, \rho_m)$ are the mixing proportions and $\varphi(x, \boldsymbol{\alpha}_{k\ell})$ is a probability density function defined on the real set \mathbb{R} .

Counts in the $r \times s$ cells of a contingency table are typically modelled as random variables. In our situation, we assume that for each block $k\ell$ the values x_{ij} are distributed according the Poisson distribution $\mathcal{P}(\alpha_i \beta_j \delta_{k\ell})$ for which the probability mass function is

$$\frac{e^{-\alpha_i \beta_j \delta_{k\ell}} (\alpha_i \beta_j \delta_{k\ell})^{x_{ij}}}{x_{ij}!}.$$

The Poisson parameter is split into α_i and β_j the effects of the row i and the column j and $\delta_{k\ell}$ the effect of the block $k\ell$. Because the aim is to maximize the complete data log-likelihood not only depending on $\boldsymbol{\theta}$ but on \mathbf{z}, \mathbf{w} , an adapted re-parametrization of the Poisson distribution becomes necessary. To this end, we impose some constraints and we assume that $\sum_{\ell} \beta_{\ell} \delta_{k\ell} = 1$ and $\sum_k \alpha_k \delta_{k\ell} = 1$ with $\alpha_k = \sum_{i,k} z_{ik} \alpha_i$ $\beta_{\ell} = \sum_{j,\ell} w_{j\ell} \beta_j$.

To tackle the simultaneous partitioning problem, we will use the CML approach, which aims to maximize the classification log-likelihood called complete data log-likelihood associated to the block mixture model. With our model, the complete data are $(\mathbf{z}, \mathbf{w}, \mathbf{x})$ and the classification log-likelihood is given by

$$L_c(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta}) = L(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z}, \mathbf{w}) = g(p(\mathbf{z}; \boldsymbol{\theta})p(\mathbf{w}; \boldsymbol{\theta})f(\mathbf{x}|\mathbf{z}, \mathbf{w}; \boldsymbol{\theta})).$$

To maximize $L_c(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta})$, like in Govaert and Nadif (2003) we propose to maximize alternatively the classification log-likelihood with \mathbf{w} and $\boldsymbol{\rho}$ fixed and

then with \mathbf{z} and $\boldsymbol{\pi}$ fixed. By noting $x_{i\ell} = \sum_j w_{j\ell} x_{ij}$, the classification log-likelihood can be written as

$$L_c(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta}) = \sum_{i,k} z_{ik} \cdot g \pi_k + \sum_{j,\ell} w_{j\ell} \cdot g \rho_\ell + \sum_{i,k} z_{ik} \sum_{\ell} x_{i\ell} \cdot g \delta_{k\ell}.$$

If we note $\mathbf{u}_i = (x_{i1}, \dots, x_{i\ell}, \dots, x_{im})$ and $\gamma_{k\ell} = x_{i\ell} \delta_{k\ell}$, the classification log-likelihood can be decomposed into two terms

$$L_c(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta}) = L_c(\mathbf{z}, \boldsymbol{\theta}/\mathbf{w}) + g(\mathbf{x}, \mathbf{w}, \boldsymbol{\rho})$$

where the first one, can be written as

$$L_c(\mathbf{z}, \boldsymbol{\theta}/\mathbf{w}) = \sum_{i,k} z_{ik} \cdot g(\pi_k \Phi(\mathbf{u}_i, \gamma_k))$$

where $\Phi(\mathbf{u}_i, \gamma_k)$ is the multinomial distribution for x_{i1}, \dots, x_{im} with the probabilities $\gamma_{k1}, \dots, \gamma_{km}$ and the second one can be written as

$$g(\mathbf{x}, \mathbf{w}, \boldsymbol{\rho}) = \sum_{j,\ell} w_{j\ell} \cdot g \rho_\ell - \sum_{\ell} x_{i\ell} \cdot g x_{i\ell}.$$

Hence, $L_c(\mathbf{z}, \boldsymbol{\theta}/\mathbf{w})$, called in the followings conditional classification log-likelihood, corresponds to the complete log-likelihood associated to a classical mixture model defined on the samples $\mathbf{u}_1, \dots, \mathbf{u}_r$. As $g(\mathbf{x}, \mathbf{w}, \boldsymbol{\rho})$ does not depend on \mathbf{z} , maximizing $L_c(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta})$ for \mathbf{w} fixed is equivalent to maximize the conditional classification log-likelihood $L_c(\mathbf{z}, \boldsymbol{\theta}/\mathbf{w})$, which can be done by the CEM algorithm applied to the multinomial mixture model. The different steps of CEM are

- E-step: compute the posterior probabilities $t_{ik}^{(c)}$;
- C-step: the k th cluster of $\mathbf{z}^{(c+1)}$ is defined with $z_{ik}^{(c+1)} = 1$ if $k = \operatorname{argmax}_{k=1, \dots, g} t_{ik}^{(c)}$ and $z_{ik}^{(c+1)} = 0$ otherwise.
- M-step: by standard calculations, one arrives at the following re-estimations parameters

$$\pi_k^{(c+1)} = \frac{\#z_k^{(c+1)}}{r} \quad \text{and} \quad \delta_{k\ell}^{(c+1)} = \frac{x_{k\ell}}{x_k \cdot x_{i\ell}}$$

where $\#$ denotes the cardinality and

$$x_{k\ell} = \sum_i z_{ik}^{(c+1)} x_{i\ell} = \sum_{ij} z_{ik}^{(c+1)} w_{j\ell} x_{ij}.$$

In the same way, we can show that

$$L_c(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta}) = L_c(\mathbf{w}, \boldsymbol{\theta}/\mathbf{z}) + g(\mathbf{x}, \mathbf{z}, \boldsymbol{\pi})$$

where

$$g(\mathbf{x}, \mathbf{z}, \boldsymbol{\pi}) = \sum_{i,k} z_{ik} \cdot g \pi_k - \sum_k x_k \cdot g x_k.$$

does not depend on \mathbf{w} and $L_c(\mathbf{w}, \boldsymbol{\theta}/\mathbf{z})$ corresponds to the complete log-likelihood associated to a classical mixture model defined on the samples $\mathbf{v}_1, \dots, \mathbf{v}_s$ where $\mathbf{v}_j = (x_{1j}, \dots, x_{kj}, \dots, x_{gj})$ with $x_{kj} = \sum_i z_{ik} x_{ij}$ and therefore develop the different steps of the CEM algorithm applied on $\mathbf{v}_1, \dots, \mathbf{v}_s$ to maximize $L_c(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta})$ for \mathbf{z} fixed.

Finally, we can describe easily the different steps of the algorithm called block CEM and noted BCEM:

1. Start from an initial position $(\mathbf{z}^{(0)}, \mathbf{w}^{(0)}, \boldsymbol{\theta}^{(0)})$.
2. Computation of $(\mathbf{z}^{(c+1)}, \mathbf{w}^{(c+1)}, \boldsymbol{\theta}^{(c+1)})$ starting from $(\mathbf{z}^{(c)}, \mathbf{w}^{(c)}, \boldsymbol{\theta}^{(c)})$:
 - (a) Computation of $\mathbf{z}^{(c+1)}, \boldsymbol{\pi}^{(c+1)}, \delta^{(c+\frac{1}{2})}$ using the CEM algorithm on the data $(\mathbf{u}_1, \dots, \mathbf{u}_r)$ starting from $\mathbf{z}^{(c)}, \boldsymbol{\pi}^{(c)}, \delta^{(c)}$.
 - (b) Computation of $\mathbf{w}^{(c+1)}, \boldsymbol{\rho}^{(c+1)}, \delta^{(c+1)}$ using the CEM algorithm on the data $(\mathbf{v}_1, \dots, \mathbf{v}_s)$ starting from $\mathbf{w}^{(c)}, \boldsymbol{\rho}^{(c)}, \delta^{(c+\frac{1}{2})}$.
3. Iterate the steps 2 until the convergence.

4 Numerical Experiments

To illustrate the behavior of our algorithms BCEM and 2CEM, we studied their performances on simulated data. We selected twenty five kinds of data arising from 3×2 -component Poisson block mixture in considering firstly the situation where the proportions are equal proportions ($\pi_1 = \pi_2 = \pi_3$ and $\rho_1 = \rho_2$). These data are obtained by varying the following parameters: the degree of overlapping which depends on the parameters $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\rho}, \delta)$, and the sizes r and s . This degree of overlapping can be measured by the Bayes error corresponding to our model. Its computation being theoretically difficult, we used Monte Carlo simulations and evaluated this error by comparing the simulated partitions and those we obtained by applying a C-step. Five degrees of overlapping have been considered and are approximatively equal to 6%, 11%, 16%, 18%, 20%. Concerning the size, we took $r \times s = (30 \times 100), (50 \times 100), (100 \times 100), (500 \times 100)$ and (1000×100) .

For each of these 25 data structures, we generated 30 samples and for each sample, we ran BCEM and CEM 100 times starting from random situations and selected the best solution for each method. In order to summarize the behavior of these algorithms, we used the proportion of misclassified points "error rate" occurring for each sample.

The results obtained are displayed in Table 1. For each data set and each algorithm, we summarize the 30 trials with the means and standard deviations of error rates obtained by comparing the partitions obtained by the both methods and the simulated partitions. In Table 2, we report the means and standard deviations of running times.

From these experiments, the main point arising are the following.

- The version 2CEM working on the two sets separately is suitably effective only when the clusters are well separated. This shows the risk of the use of such methods when the clusters are ill-separated.

Table 1. Comparison of BCEM and 2CEM for 30 kinds of data : means and standard deviations of error rates

Size	Overlap					
	1	2	3	4	5	
30	BCEM	0.177 (0.084)	0.321 (0.186)	0.560 (0.164)	0.665 (0.106)	0.657 (0.135)
	2CEM	0.309 (0.066)	0.427 (0.134)	0.625 (0.124)	0.663 (0.092)	0.678 (0.101)
50	BCEM	0.105 (0.055)	0.239 (0.076)	0.488 (0.126)	0.707 (0.116)	0.682 (0.146)
	2CEM	0.262 (0.066)	0.350 (0.090)	0.581 (0.103)	0.701 (0.086)	0.710 (0.102)
100	BCEM	0.063 (0.024)	0.155 (0.015)	0.335 (0.062)	0.449 (0.160)	0.623 (0.155)
	2CEM	0.183 (0.056)	0.281 (0.049)	0.477 (0.101)	0.570 (0.086)	0.658 (0.124)
500	BCEM	0.061 (0.011)	0.123 (0.011)	0.166 (0.019)	0.198 (0.022)	0.255 (0.040)
	2CEM	0.098 (0.019)	0.195 (0.024)	0.277 (0.043)	0.375 (0.070)	0.446 (0.080)
1000	BCEM	0.065 (0.005)	0.118 (0.007)	0.162 (0.012)	0.187 (0.016)	0.212 (0.029)
	2CEM	0.083 (0.012)	0.190 (0.022)	0.247 (0.025)	0.300 (0.052)	0.376 (0.037)

Table 2. Comparison of BCEM and 2CEM for 30 kinds of data : means and standard deviations of running times

Size	Overlap					
	1	2	3	4	5	
30	BCEM	2.102 (0.126)	2.162 (0.187)	1.934 (0.176)	1.870 (0.131)	1.871 (0.094)
	2CEM	1.422 (0.058)	1.490 (0.048)	1.565 (0.198)	1.476 (0.039)	1.461 (0.037)
50	BCEM	2.314 (0.274)	2.901 (0.173)	2.823 (0.553)	2.394 (0.098)	2.444 (0.084)
	2CEM	2.440 (0.150)	2.418 (0.141)	2.689 (0.693)	2.437 (0.149)	2.349 (0.157)
100	BCEM	2.282 (0.147)	3.386 (0.192)	3.785 (0.335)	3.230 (0.225)	2.827 (0.169)
	2CEM	4.599 (0.071)	4.607 (0.070)	4.685 (0.061)	4.653 (0.104)	4.560 (0.053)
500	BCEM	6.346 (0.435)	7.387 (0.758)	8.784 (0.933)	7.800 (0.833)	6.868 (0.729)
	2CEM	26.760 (0.250)	26.430 (0.227)	26.719 (0.364)	26.540 (0.436)	26.407 (0.183)
1000	BCEM	9.460 (1.130)	10.521 (0.981)	10.189 (0.874)	8.382 (0.609)	7.916 (0.626)
	2CEM	54.566 (0.280)	54.453 (0.318)	54.387 (0.348)	54.796 (0.443)	54.277 (0.186)

- Incontestably BCEM outperforms 2CEM. The results are very encouraging and its performance increases with the size of data.
- It appears clearly that BCEM is undoubtedly faster as soon as the size is large enough.

We carried out other simulations on large data sets with proportions dramatically different, not included in this text, which confirms these remarks.

5 Conclusion

Setting the problem of block clustering under the CML approach, we have compared block CEM and two-way CEM. The first one gives encouraging results on simulated data and real data and is therefore strongly recommended : it is faster and better than two-way CEM. Currently, we are evaluating block CEM on other large real data sets. In this paper, we have considered the block clustering for contingency tables under the CML approach and, as in Govaert and Nadif (2005a, 2005b) for binary data, it would be interesting to study the block clustering of contingency table under the ML and fuzzy approaches.

References

- Arabie, P., J., H.L.: The bond energy algorithm revisited. *IEEE Transactions on Systems, Man, and Cybernetics* **20** (1990) 268–274
- Bock, H.: Simultaneous clustering of objects and variables. In Diday, E., ed.: *Analyse des Données et Informatique*, INRIA (1979) 187–203
- Celeux, G., Govaert, G.: A classification em algorithm for clustering and two stochastic versions. *Computational Statistics and Data Analysis* **14** (1992) 315–332
- Cheng, Y., Church, G.: Biclustering of expression data. In: *Proceedings of the Eighth International Conference on Intelligent Systems for Molecular Biology (ISMB)*. (2000) 93–103
- Dempster, A.P., Laird, N.M., Rubin, D.B.: Maximum likelihood from incomplete data via the em algorithm (with discussion). *Journal of the Royal Statistical Society* **B 39** (1977) 1–38
- Dhillon, I.: Co-clustering documents and words using bipartite spectral graph partitioning. In: *ACM SIGKDD Conference, San Francisco, USA*. (2001) 269–274
- Duffy, D.E., Quiroz, A.J.: A permutation-based algorithm for block clustering. *Journal of Classification* **8** (1991) 65–91
- Garcia, H., Proth, J.M.: A new cross-decomposition algorithm: The GPM comparison with the bond energy method. *Control and Cybernetics* **15** (1986) 155–165
- Govaert, G.: *Classification croisée*. Thèse d'état, Université Paris 6, France (1983)
- Govaert, G.: Simultaneous clustering of rows and columns. *Control and Cybernetics* **24** (1995) 437–458
- Govaert, G., Nadif, M.: Clustering with block mixture models. *Pattern Recognition* **36** (2003) 463–473
- Govaert, G., Nadif, M.: An EM algorithm for the block mixture model. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **27** (2005) 643–647
- Govaert, G., Nadif, M.: Fuzzy clustering to estimate the parameters of block mixture models. *Soft Computing* (in press, 2005)
- Hartigan, J.A.: *Clustering Algorithms*. Wiley, New York (1975)
- Marchotorchino, F.: Block seriation problems: A unified approach. *Applied Stochastic Models and Data Analysis* **3** (1987) 73–91
- McLachlan, G.J., Peel, D.: *Finite Mixture Models*. Wiley, New York (2000)
- Ritschard, G. Zighed, D., Nicoloyannis, N., Maximisation de l'association par regroupement de lignes ou de colonnes d'un tableau croisé. *Revue de Mathématiques & Sciences Humaines* **39** (2001) 81–97
- Symons, M.J.: Clustering criteria and multivariate normal mixture. *Biometrics* **37** (1981) 35–43

An Imbalanced Data Rule Learner

Ca. h Ha. Ngu e. a. d Tu Ba. H.

School of Knowledge Science,
Japan Advanced Institute of Science and Technology,
1-1 Asahidai, Nomi, Ishikawa, 923-1292, Japan
{canhhao, bao}@jaist.ac.jp

Abstract. Imbalanced data learning has recently begun to receive much attention from research and industrial communities as traditional machine learners no longer give satisfactory results. Solutions to the problem generally attempt to adapt standard learners to the imbalanced data setting. Basically, higher weights are assigned to small class examples to avoid their being overshadowed by the large class ones. The difficulty determining a reasonable weight for each example remains. In this work, we propose a scheme to weight examples of the small class based solely on local data distributions. The approach is for categorical data, and a rule learning algorithm is constructed taking the weighting scheme into account. Empirical evaluations prove the advantages of this approach.

1 Introduction

In a practical sense, applying standard machine learning methods to real world tasks where their class distribution is imbalanced is problematic. This is the case of a data set, in which the number of false positives of the class is substantially smaller than the others. For instance, if the false class accounts for only 2% of the total number of examples in the data set, a classifier can get a high accuracy of 98% just by assigning a large weight to the large class. However, in such a case, the classifier completely fails to learn the small class, which is usually of interest. In practice, researchers have encountered this problem in many applications, including the detection of fraudulent transactions [1], network intrusion detection [2] and missile satellite radar images [3].

The reason that standard classifiers cannot learn a satisfactory performance on such data sets is because they make the fundamental assumption that frequencies of classes are equally distributed. Adaptations to imbalanced data sets are usually made by giving small class examples higher weights. One simple approach is resampling, which duplicates small class examples to form a subset of large class ones. Such approaches do not have high performance, as reported in [4], because larger examples are affected differently by the class imbalance problem. SMOTE [5] combines synthetic example generation with down-sampling, but the resampling degree is not specified. Resampling to reflect relative weights between classes still remains a art.

It is believed that better preprocessing to weight examples differently, and various approaches have been proposed. Kubat et al. [3] insist that large class examples

1. a mini ed regi . sh u d be e ighte d zero as . g as that 1 creases perf rma ce measure. Lear 1 g . a custer basis is used [6] t e ight e amp es acc rdi g . Lear 1 g decis i trees (DT) [7] is made 1 depe de t f c ass freque cies b usi g the Area Under the ROC Curve (AUC) as a sp itti g criteri , hich is equi a e t t e ighti g e amp es acc rdi g t their distributi i the set f e amp es c e red b the sp itti g . des. A ge era a t . ptima e ight e amp es (1 Ba es risk se se) is usi g MetaC st [8], b baggi g a d the pr b abilit estimati . H e er, 1 high imba a ced data sets, e amp es f sma c ass are rare ear ed, mak i g their ptima c sts e treme high. Agai , it is sti a cha e ge t e ight e amp es ptima f r the imba a ced data pr b em.

We pr p se a meth d t estimate the ptima e ight f each sma c ass e amp e basi g s e . ca data distributi s. The i tuti is that b . ki g m re c se i t . ca data distributi , e ha e m re cha ce t re ea usefu i f rmati ab ut the e ect f c ass imba a ce. T this e d, e first defi e the c ept f . . . , hich characterizes . ca data distributi a d the determi es e amp es' eights ith the aim g f ma imizi g AUC i the i ci t . The e ight is i tegrated i t a ru e i ducti a g rithm at the ru e pr i g step.

The paper is rga ized as f . s. Secti . 2 is the f u dati a d f rmu ati . f ur . ca adapti e e ighti g scheme. I tegrati . f the e ighti g scheme i t . ur ru e ear 1 g a g rithm is described i Secti . 3. I Secti . 4, e sh e perime ta e a uati s f the scheme t . ther imba a ced data c assifiers. C . c usi s are prese ted a d future . rk is discussed i secti . 5.

2 Locally Adaptive Weighting Scheme

We appr ach the imba a ced data pr b em b gi i g a e ight adapti e f r each sma c ass e amp e, hie keepi g the eights f arge c ass e amp es at a defau t a ue (i.e. 1). The ke idea is t e ight each sma c ass based . its . ca . eighb rh d (he ce, it is . ca adapti e), hich is defi ed as the i ci t f the ru e c e ri g it. This secti . 1 defi e the c ept f i ci t a d deri e the f rmu ati . f e amp e e ighti g based . i ci t usi g AUC as the criteri .

Vicinity: The idea behi d i ci t is as f . s. C . sider t . ru es $R_i, i = 1, 2$ ith the same c e rage f r e r c ass (R_i c e rs n_i, p_i e amp es fr m arge a d sma c asses respecti e , $n_1 = n_2, p_1 = p_2$). C . e t i a , the t . ru es are e a uated as the same g d ess (e.g., precisi f r sma c ass $\frac{p_i}{p_i+n_i}$). Assume that e ha e s me a t defi e the surr u d f a ru e, ca ed a eighb rh d. If R_1 is ike t be pr u ed t a better . e, the i its eighb rh d there must be s me e amp es f the same c ass as the predicti g c ass f R_1 . O the ther ha d, R_2 is surr u ded b e amp es fr m ther c asses, he ce it ca . t be pr u ed t a better . e. Our idea is t e a uate the t . ru es di ere t , R_1 t be higher tha R_2 , re ecti g their abilit t be pr u ed. This di ere t e a uati . is based . the fact that there is a set f e amp es i each ru e's eighb rh d, hich creates the di ere ce i pr u i g abilit . B i ci t , the , e mea this set f e amp es. We defi e i ci t based . the c ept f k- i ci t .

Definition: For a rule R and a data set D , the k -neighborhood of R is defined as:

Definition: For a rule R and a data set D , the k -vicinity of R is defined as:

$$vicinity(R, k) = \{x \mid x \in D, Distance(R, x) \leq k\} \tag{1}$$

k -vicinity is a subset the training data set, which is partially covered by the rule after k steps of generalization. The smaller k is, the higher the coverage the examples k -vicinity makes in the generalization (pruning) ability of the rule. For example, 0-vicinity is the set of examples covered by the rule, k -vicinity is the whole data set if m is the number of attribute-value pairs in the rule body. The set of k -vicinities is a nested chain of subsets of the data set, meaning: $vicinity(R, 0) \subseteq vicinity(R, 1) \subseteq \dots \subseteq vicinity(R, m)$. We define k -vicinity using this chain with heights. Formally, k -vicinity is a function f over k -vicinities.

$$vicinity = f\{vicinity(R, 0), vicinity(R, 1), \dots, vicinity(R, m)\} \tag{2}$$

Estimating k -vicinity is a difficult task. However, the standard problem is to let the k -vicinity of a rule remain a virtual concept. We need to calculate the k -vicinity of a rule R as the height g scheme discussed in the next section.

Example Weighting and Rule Evaluation: As a k -vicinity is expected to contain examples that achieve the pruning ability of a rule, we use this assumption to define the k -vicinity of a rule as the best g -optimal assignment with its k -vicinity. Our idea is to weight examples in the k -vicinity such that the optimal assignment coincides with the best assignment cost. Defining optimal assignment as a k -vicinity results in a local adaptive height g scheme.

We define optimal assignment as the one that gives the largest AUC. AUC is a popular metric for comparing classifiers' performance [9, 10] where the misclassification costs are unknown. When a classifier is a set of rules, as in Figure 1 (a), the ROC curve contains a set of line segments. Here, the classifier is assumed to have features, sorted in decreasing order of their precision for a class. For simplicity, we assume that there is only one rule for small class k -vicinity. Then, the ROC curve of a classifier (between R_1 and R_2) in its k -vicinity would look like Figure 1 (b). The classifier here consists of a rule (say R) and the default rule predicting the large class. Suppose R covers p small and n large class examples, and the k -vicinity contains P small and N large class examples. The rule evaluation metric, defined to be AUC above, is calculated [7] as:

$$AUC(R) = \frac{p}{2P} - \frac{n}{2N} + \frac{1}{2} \tag{3}$$

The above formula implies that the height of a small class example in this k -vicinity is $\frac{N}{P}$ where the height of a large class example is the default value 1.

The rule evaluation metric is used to compare different rules for search bias. However, it is not natural to compare AUC in different contexts (vicinities).

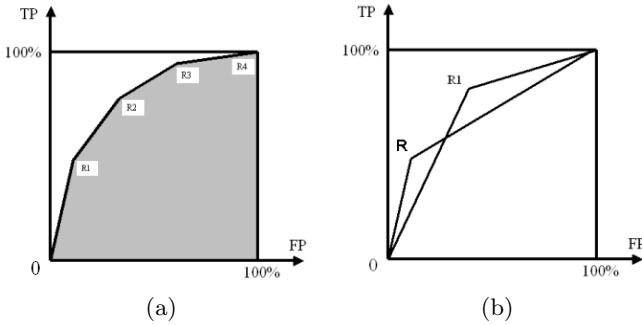


Fig. 1. The ROC space: (a) plot of a rule learner. (b) for rule comparison in a vicinity.

Here, we propose a comparison strategy that rule R_1 is considered better than rule R if and only if it gives a higher AUC than the one of R . Equivalently, R_1 is considered better than R if their AUC difference (formula 4) is positive.

$$AUC(R_1) - AUC(R) = \frac{1}{2} \left(\frac{p_1 - p}{P} - \frac{n_1 - n}{N} \right) = \frac{1}{2} \left[(p_1 - p) \frac{N}{P} - (n_1 - n) \right] \frac{1}{N} \quad (4)$$

For the purpose of comparing rules for search bias, it is sufficient to know $\frac{N}{P}$. This is the reason we are interested to remain an abstract concept, here we heuristically estimate $\frac{N}{P}$ directly. From equation 2, we propose to estimate class distribution ratio $\frac{N}{P}$ in the vicinity to be:

$$\frac{N}{P} = \sum_{k=1}^m w_k * \frac{N_k}{P_k} \quad (5)$$

where m is the number of classes. In this formula, $\frac{N_k}{P_k}$ is class distribution ratio in k -class and w_k is its associated weight, $\sum w_k = 1$. This just means that the class distribution ratios of different classes to estimate that of the vicinity, in similar fashion to a shrinkage estimator, to make it robust. The definition of vicinity is two adjacent neighborhood scheme, namely the set $\{w_k\}$. If w_k is large for small k s, vicinity reflects more carefully from it. If w_k is large for large k s, vicinity is more global. If we take the definition to be the whole data set, then $w_m = 1$. Having taken the set $\{w_k\}$ is a generalization of a simple constant sensitivity classification. In this algorithm, we fix the default as:

$$w_k = \frac{1}{m}, k = \overline{1, m} \quad (6)$$

Such weights can also be set adaptively by users.

Discussion: The key point which makes this simple neighborhood scheme suitable for imbalanced data is the use of a local neighborhood. Having a macro view around a rule gives us a better picture of how much the imbalance may hinder classification rules. Empires far away from the boundaries of classes may not participate in a vicinity, as a result affecting classification (this is similar to the idea behind SVMs).

IDL

1. Generate a candidate rule set
2. Prune rules from high coverage to low

GenerateRuleSet

1. Generate a decision tree
2. Stop when leaf nodes contain only example from one class
3. Extract the set of leaf nodes that contain only examples of small class
4. Convert those nodes into rules and return

PruneRules

1. Sort rules according to coverage
2. From high to low coverage rule do
 3. Remove *best* attribute value pair
 4. Until no more AUC is gained
5. Return pruned rules

Fig. 2. IDL algorithm

3 IDL: Imbalanced Data Learner

IDL is a rule inductive algorithm, which uses a one-sided selection strategy, taking example weights into account. It learns rules for a small class. IDL consists of two steps. First, it generates a set of candidate rules for the small class by growing a decision tree, which are meant to be complete and of high precision. Then it prunes these rules by greedily removing attribute value pairs that make them robust. Example weighting is used in rule pruning step using Formula 3 as the rule evaluation metric. The overall strategy is depicted in Figure 2.

In the candidate rule set generation step, IDL grows a decision tree and stops when the leaf nodes contain examples from one class. As recommended in [11], IDL takes the impurity $(2\sqrt{p(1-p)})$ gain as the splitting criterion. After the decision tree is fully grown, the set of leaf nodes that contain only small class examples are collected and turned into a set of rules. In the second step, the collected rules for the small class are sorted in decreasing order of coverage. Starting from the highest coverage rule, each rule is pruned by removing the best attribute value pair (the one having the highest AUC decrease), according to Formula 4. It stops when removing does not improve either AUC or the precision of the rule (calculated with taking weights into account) falls under a certain threshold. The examples covered by a rule are marked so that they are covered again, so that their weights are retained. This makes the rules overlap, and as a great improvement the recall of the classifier. The threshold represents the minimum precision a rule should achieve, reflecting the amount of noise in the data. This threshold is generated separately by the users, and is estimated in IDL as follows. First, set it to 80%, then do a 10-fold stratified cross-validation.

... the data set to estimate its difficult to learn. Taking the F-measure as the small class, say f (percentage), then the threshold takes the value $max(50, f - 10)$.

In the first step, IDL constructs a unpruned decision tree, which is of $O(ea)$ time complexity, where e is the number of examples and a is the number of attributes. In the second step, support it generates k rules, each has maximum n_k attribute value pairs. As each pruning operation requires a pass of the data set to calculate the class distribution ratio in the vicinity, the time complexity of this step is at most $O(ekn_k)$.

4 Experimental Evaluation

We evaluated IDL in its ability to learn a small class, and compared it to other approaches. The first of these was SMOTE-NC (see C4.5) [5], the most accurate version of what is arguably best method (SMOTE) for learning imbalanced data. Since SMOTE is sensitive to its degree of sampling parameters, we ran it in three degrees of small class up-sampling, namely $N=100\%$, $N=300\%$ and $N=700\%$. We also compared IDL to a general classifier of C4.5, with and without cost sensitive settings. In cost sensitive settings (C.S.), the relative cost is just the ratio of class distribution of the data set. Building is a successful feature handling imbalanced data learners [12], such as AdaBoost over C4.5 was also compared. We used these classifiers from WEKA¹. All algorithms ran with their default parameters. We used F-measure as the small class as the performance criterion, instead of the AUC measure (since we are only using small class rules).

$$F - measure = \frac{2 * precision * recall}{precision + recall} \tag{7}$$

We evaluated these algorithms on selected fifteen UCI data sets², where small class as chosen to be small class, and the other classes were merged to become the large class. As the algorithm is for categorical data, all data sets were discretized. We split data sets with a ratio of 75-25 random into a stratified manner, the large parts were used for training and the small parts for testing. Table 1 shows the results of testing on the small part of the data. The columns are names, percentage of small class preceded with class index and the classifiers (SMOTE is tested with three parameters). All numbers are in percentage. The last line shows average performance for all data sets.

The table shows that our approach outperforms general classifiers by a large margin, and is competitive to a three parameter for SMOTE. IDL shows an improvement of 11.74% in terms of F-measure as the small class compared to a standard classifier of C4.5. For the cost sensitive settings of C4.5 (C.S.), it also improves by 3.81%. Compared to AdaBoost, IDL's accuracy is 2.85% higher. This means that IDL is more suitable for imbalanced data than general classifiers. Comparing IDL to imbalanced data learner of SMOTE (SMOTE-NC version), IDL is a competitive to the three parameter settings; the average

¹ www.cs.waikato.ac.nz/ml/weka/
² <http://www.ics.uci.edu/mllearn/MLRepository.html>

Table 1. Comparison of Classifiers on UCI data

Name	%	C4.5	C.S.	SMOTE				A.Boost	IDL
				100	300	700	average		
annealing1	11.0	73.9	62.3	77.4	71.0	66.7	71.7	70.6	96.2
car3	3.7	66.7	84.2	77.4	80.0	80.0	79.1	80.0	76.9
flare4	8.0	0.0	36.9	30.4	36.1	38.6	35.0	32.7	29.0
glass3	13.5	93.3	73.7	93.3	82.4	82.4	86.0	80.0	85.7
hypo0	5.0	85.7	81.9	85.7	83.1	83.1	84.0	84.6	84.6
inf0	6.3	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
krkopt16	0.9	85.7	82.8	88.4	88.4	90.1	89.0	84.1	87.1
krkopt4	0.7	58.0	69.1	66.7	71.8	66.7	68.4	61.3	75.9
led7	8.4	59.8	59.9	63.9	61.6	50.5	58.7	50.7	62.4
letter0	3.9	91.0	90.0	92.0	91.8	89.7	91.2	96.0	92.0
satimage3	9.7	51.5	52.8	57.9	51.8	51.3	53.7	57.9	50.3
segmentation5	14.1	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
sick1	6.5	82.8	69.6	82.8	81.8	73.8	79.5	82.8	80.9
vowel5	9.1	0.0	75.2	67.8	69.4	65.5	68.2	80.0	60.0
yeast4	3.4	0.0	30.8	33.3	44.4	40.0	39.2	21.1	43.5
Average	6.94	63.23	71.16	74.59	74.24	71.89	73.57	72.12	74.97

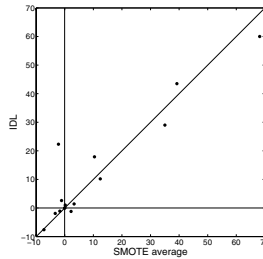


Fig. 3. Improvement of IDL versus SMOTE

performance of SMOTE is 1.40% better than that of IDL. It is interesting that there is no systematic way to determine the resampling degree for SMOTE.

It is interesting to look at the improvement of IDL over C4.5 compared to the average improvement of that of SMOTE over C4.5 in Figure 3. X-axis is the performance improvement of SMOTE (average parameters), while y-axis is for IDL. The set of points shows a clear linear relationship. This means that improvement of IDL is proportional to that of SMOTE, meaning that IDL is consistently similar to SMOTE.

5 Conclusion

We have proposed a method to integrate ensemble for a small class based on their local neighborhood. Neighborhood is defined as the virtual concept of locality, where computation is based on k-nearest neighbors. The algorithm is clear and

accurate than general classifiers, including AdaBoost and MetaC. SMOTE also has the advantage of not requiring resampling parameters. From this, we can conclude that the 1-f neighbor 1-f neighbor distance sampling is useful for handling imbalanced data.

The clear limitation of this method is how to define the neighbor scheme for a class. For the moment, computation is its main problem, which should be reduced for large data sets. Applying the neighbor scheme to other classifiers for imbalanced data is a natural extension. Whether this data distribution can be used to improve classifiers in general is an open question.

References

- [1] Fawcett, T., Provost, F.: Combining data mining and machine learning for effective user profiling. In Simoudis, Han, Fayyad, eds.: The Second International Conference on Knowledge Discovery and Data Mining, AAAI Press (1996) 8–13
- [2] Lazarevic, A., Ertöz, L., Ozgur, A., Srivastava, J., Kumar, V.: "evaluation of outlier detection schemes for detecting network intrusions". In: Third SIAM International Conference on Data Mining. (2003)
- [3] Kubat, M., Holte, R.C., Matwin, S.: Machine learning for the detection of oil spills in satellite radar images. *Machine Learning* **30** (1998) 195–215
- [4] Japkowicz, N.: The class imbalance problems: Significance and strategies. In: Proceedings of the 2000 International Conference on Artificial Intelligence (IC-AI'2000). Volume 1. (2000) 111–117
- [5] Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P.: Smote: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research* **16** (2002) 321–357
- [6] Nickerson, A., Japkowicz, N., Milios, E.: Using unsupervised learning to guide re-sampling in imbalanced data sets. In: Eighth International Workshop on AI and Statistics. (2001) 261–265
- [7] Ferri, C., Flach, P., Hernandez-Orallo, J.: Learning decision trees using the area under the roc curve. In Sammut, C., ed.: Nineteenth International Conference on Machine Learning ICML'02, Morgan Kaufmann (2002)
- [8] Domingos, P.: Metacost: A general method for making classifiers cost-sensitive. In: Knowledge Discovery and Data Mining. (1999) 155–164
- [9] Provost, F., Fawcett, T.: Robust classification for imprecise environments. *Machine Learning* **42** (2001) 203–231
- [10] Furnkranz, J., Flash, P.: An analysis of rule evaluation metrics. In: The Twentieth International Conference on Machine Learning (ICML'03), AAAI Press (2003) 202–209
- [11] Elkan, C.: The foundations of cost-sensitive learning. In: Seventeenth International Joint Conference on Artificial Intelligence (IJCAI'01). (2001) 973–978
- [12] Joshi, M.V., Agarwal, R.C., Kumar, V.: Predicting rare classes: can boosting make any weak learner strong? In: Proceedings of the eighth ACM international conference on Knowledge discovery and data mining, ACM Press (2002) 297–306

Improvements in the Data Partitioning Approach for Frequent Itemsets Mining

Son N. Nguyen and Maria E. Orlowska

School of Information Technology and Electrical Engineering,
The University of Queensland, QLD 4072, Australia
{nson, maria}@itee.uq.edu.au

Abstract. Frequent Itemsets mining is well explored for various data types, and its computational complexity is well understood. There are methods to deal effectively with computational problems. This paper shows another approach to further performance enhancements of frequent items sets computation.

We have made a series of observations that led us to inventing data pre-processing methods such that the final step of the Partition algorithm, where a combination of all local candidate sets must be processed, is executed on substantially smaller input data. The paper shows results from several experiments that confirmed our general and formally presented observations.

Keywords: Association rules, Frequent itemset, Partition, Performance.

1 Introduction

Since the association rules mining introduction by Argawal et al. [5], many algorithms and their subsequent improvements have been proposed to solve association rules mining, especially frequent itemsets mining problems.

In this paper, we review the state of the art in association rules mining with a focus on frequent itemsets mining. There are many well-accepted approaches such as “Apriori” by Argawal et al. [1], ECLAT by Zaki [7], and more recently “FP-growth” by Han et al. [8]. Another interesting class of solutions is based on the data partitioning approach. This fundamental concept was originally proposed as a Partition algorithm by Savaserse et al. [2], and it was improved later in AS-CPA by Lin et al. [4] and ARMOR by Pudi et al. [11]. A common feature of these results is their target, namely the limitation of I/O operations by considering data subsets dictated by the main memory size.

An intriguing question is whether we could improve the overall performance of mining large data sets by a smarter but not too ‘expensive’ design of the data fragments - rather than determine them by a sequential transaction allocation based on the fragment size only.

The main goal of this paper is to demonstrate our observations, generalize, and specify corresponding data pre-processing for the Partitioning approach in order to improve the performance. Our study is supported by a series of experiments which indicate a dramatic improvement in the performance of the Partitioning approach with our fragmentation method, in contrast to the traditional one [2].

The remainder of the paper is organised as follows. Section 2 introduces the basic concepts related to frequent itemsets mining. Section 3 reviews the current state of art in the field, especially for frequent itemsets mining and the Partitioning approach. Section 4 presents our observations and open issues. We propose the pre-processing data fragmentation solution in section 5. Section 6 shows the result from our experiment, and finally, we present our concluding remarks in section 7.

2 Preliminary Concepts

For the completeness of this presentation and to establish our notation, this section gives a formal description of the problem of mining frequent itemsets. It can be stated as follows [1]:

Let $I = \{i_1, i_2, \dots, i_m\}$ be a set of m distinct literals called items. Transaction database D is a set of variable length transactions over I .

Each transaction contains a set of items $\{i_j, i_k, \dots, i_h\} \subseteq I, i_j < i_k < \dots < i_h$. Each transaction has an associated unique identifier called TID.

For an itemset $X \subseteq I$, the support is denoted $\text{sup}_D(X)$, equals to the fraction of transactions in D containing X .

The problem of mining frequent itemsets is to generate all frequent itemsets X that have $\text{sup}_D(X)$ no less than user specified minimum support threshold.

3 Review Frequent Itemsets Mining

Throughout the last decade, there have been many attempts and well-known algorithms that target an efficient solution of the frequent itemsets mining problem. However, the performance of these algorithms depends on many, often very specific input data features and additionally, implementation environments. As a result, several claims made in earlier papers were later debated by other authors.

3.1 Partitioning Approach for Frequent Itemsets Mining

Savasere et al. [2] proposed the Partition algorithm based on the following principle. A fragment $P \subseteq D$ of the database is defined as any subset of the transactions contained in the database D . Further, any two different fragments are non-overlapping. *Local support* for an itemset is the fraction of transactions containing that itemset in a fragment. *Local candidate itemset* is being tested for minimum support within a given fragment. A *Local frequent itemset* is an itemset whose local support in the fragment is no less than the minimum support. *Global support, Global candidate itemset, Global frequent itemset* are defined as above except they are in the context of the entire database. The goal is to find all *Global frequent itemsets*.

The following Lemma 1 supports the main principle of the Partition algorithm.

Lemma 1: If X is a frequent itemset in database D , which is partitioned into n fragments P_1, P_2, \dots, P_n , then X must be a frequent itemset in at least one of the n fragments.

Proof: Due to the limit space, the proof can be seen in [10]

The Partition algorithm divides D into n fragments. The algorithm first scans fragment P_i in the main memory at a time, for $i = 1, \dots, n$, to find the set of all *Local frequent itemsets* in P_i , denoted as LP_i . Then, by taking the union of LP_i , a set of candidate itemsets over D is constructed, denoted as C^G . Based on *Lemma 1*, C^G is a superset of the set of all *Global frequent itemsets* in D . Finally, the algorithm scans each fragment for the second time to calculate the support of each itemset in C^G and to find the *Global frequent itemsets*.

3.2 Related Work in Partitioning Approach

One of the Partition algorithm derivatives is AS-CPA (*Anti-Skew Counting Partition Algorithm*) by Lin et al. [4]. Recently, there has been another development based on the partitioning approach in the ARMOR algorithm by Pudi et al. [11].

All the above algorithms mainly attempt to reduce the number of false candidates as early as possible. However, they do not consider any features and characteristics of data sets in order to partition the original data set more suitably for further processing.

Further in this paper, we demonstrate that looking more closely into the data itself may deliver good gains in overall performance. As a result, the *Local frequent itemsets* can be dramatically reduced. Furthermore, in many cases that leads to a larger number of common Global candidates among fragments. Finally, as a consequence, this approach reduces substantially the *Global candidates* (C^G set).

4 Observations in the Partitioning Approach

We begin by considering the first and very obvious measurable data-partitioning attribute – the size of fragments and their impact on the efficiency of the frequent items search process. Further on we examine more closely the composition of fragments at the design time to ensure that selection of transactions satisfy some desired properties.

4.1 Reasoning About Size of Fragment

It is not hard to observe that the size of the fragments is inverse-proportional to the size of the output of Local computation. Hence, the question is: *What is a ‘good’ fragment size?* We consider several heuristic methods to identify the suitable size of fragments.

We note the following observation: the smaller fragment generates a more negative effect on the number of *Local frequent itemsets*. Clearly, the best partitioning of data set D into n fragments is defined as a method that generates the smallest number of Global candidates. We denote this smallest number as G_n . Note that the perfect solution would have to exhibit the following property; *every fragment of the data generates identical Local frequent itemsets*.

We generalise these observations as follows;

Lemma 2: If database D is partitioned into $(n+1)$ fragments P_1, P_2, \dots, P_{n+1} then the number of Global candidates, denoted $|C_{n+1}^G|$, is always greater than or equal to G_n ; $|C_{n+1}^G| \geq G_n$

Proof: Due to the limit space, the proof can be seen in [10]

As a consequence, the size of a fragment should maintain proper balance in order to control the number of *Local frequent itemsets*.

4.2 Some Characteristics of Fragment Data

Data skew has a negative impact on the Partitioning approach. Basically, data skew causes the *Local frequent itemsets* generated from different fragments to have very few common elements. In such situations, the number of *Global candidates* (being the union of all LP_i) is rather large.

Obviously, fragments that have many dissimilar transactions (transactions with small or empty intersections) generate a small number of *Local frequent itemsets*. In this paper we call them *dissimilar fragments*.

These observations confirm our initial hypothesis that there are some relationships between the composition of fragments and the amount of computation required at the end. We illustrate the fact that a larger number of fragments increase the size of the computation space. In addition, for given number of fragments n , a different partition also impacts on the number of *Global candidates*. Furthermore, the gap in performance is increased dramatically when the support threshold is decreased and the number of fragments is increased.

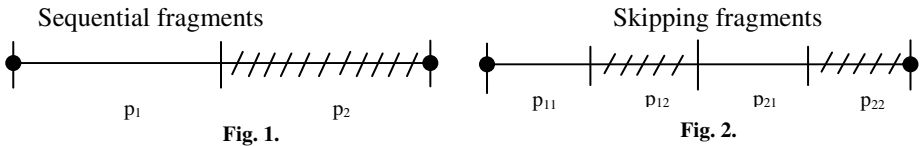
5 Data Set Pre-processing

We present the following algorithms for original data pre-processing.

5.1 Naive Algorithm

One of the simplest techniques to be considered is the skipping technique. Before formalising this concept we show a simple example to illustrate its main principle.

Consider data set D represented by a straight line on figures below. We partition D into 2 fragments as illustrated on the Figure 1. When a Skipping technique is used then D is initially divided into 4 small sequential parts. Each fragment is created by taking the union of 2 small skipping parts as it is shown on Figure 2.



One can easily generalise such partition process for any higher number of expected fragments.

5.2 An incremental Clustering Algorithm

The incremental clustering algorithm is our idea for pre-processing. The data set will be scanned only once and all clusters (fragments) containing mostly dissimilar transactions are generated at the end of that scan. We introduce some basic definitions.

Definition 5.1: Cluster Centroid is a set of all Items in the cluster, we denote it $C_i = \{I_1, I_2, \dots, I_n\}$. Additionally, each item in C_i has its associated weight which is its number of occurrences in the cluster; $\{w_1, w_2, \dots, w_n\}$

Definition 5.2: Similarity function between two item sets, in particular a transaction and a cluster centroid, is denoted $\text{Sim}(T_i, C_j)$ and defined as follows;

$\text{Sim}(T_i, C_j) \rightarrow \mathbb{R}^+$; Calculation of this function:

1. Let S be the intersection between the arguments of Sim function, $S = T_i \cap C_j$
2. If $S = \emptyset$ then $\text{Sim}(T_i, C_j) = 0$. Otherwise, $S = \{I_1, I_2, \dots, I_m\}$ with the corresponding weights $\{w_1, w_2, \dots, w_m\}$ in cluster C_j , respectively, therefore $\text{Sim}(T_i, C_j) = w_1 + w_2 + \dots + w_m$

Cluster Construction:

Informally, each transaction is evaluated in terms of the following criteria;

- a) We assign a new transaction T_i to cluster C_j which has the minimum $\text{Sim}(T_i, C_j)$ value among open clusters (*a cluster is open if has not exceeded its expected size in terms of number of transactions*).
- b) Each new allocation to a cluster C_j , updates the cluster centroid C_j . All already existing common items' weight is increased by 1, and the other new items are added to C_j with the weight of value 1.

Reasoning about the size of clusters: based on the observation in section 4.1, the cluster sizes should be well balanced.

The pseudo incremental clustering algorithm is described as the following

Input: Transaction database: D ; k – number of output clusters

Output: k clusters based on the above criteria for Partition approach.

Begin

1. Assign the first k transactions to all k clusters, and initialize the all Cluster Centroids: $\{C_1, C_2, \dots, C_k\}$
2. Consider the next k transactions: $\{T_1, T_2, \dots, T_k\}$. These k transactions are assigned to k different clusters. These operations are done based on the following criteria: (i) the minimum similarity between the new transaction and the suitable clusters; (ii) the sizes of these clusters are controlled to keep the balance. The following are more detail about this processing.

Let $C^{\text{run}} = \{C_1, C_2, \dots, C_k\}$ is a set of all k clusters; $T^{\text{run}} = \{T_1, T_2, \dots, T_k\}$

For each transaction T_i in T^{run} : T_1 to T_k

Begin

- a) Calculate the similar functions between T_i and all the clusters in C^{run} ; determine the minimum similar function value, denoted $\text{Sim}(T_i, C_j)$
- b) Assign T_i to cluster C_j which has the minimum $\text{Sim}(T_i, C_j)$ value. Update the cluster centroid C_j
- c) Remove C_j from the set of all the suitable clusters in order to keep the same size constraint. $C^{\text{run}} = C^{\text{run}} - \{C_j\}$;

End

3. Repeat step 2 till all transactions in D are clustered

End

The time complexity of this incremental clustering algorithm is about $O(|D| * k * m)$ where $|D|$ is the number of all transactions, k is the given number of clusters, and m is the number of all items in D .

6 Experiments

In this section, we conducted experiments on: one synthetic data set [1], and 3 real data sets [13]. These data sets are converted to format as the above definitions.

Table 1. The characteristics of data sets

Data sets	Transactions	Items	DB Size (~MB)
T10I4D100K	100K	870	4
WebView-1	26K	492	0.7
WebView-2	52K	3335	2
BMS-POS	435K	1657	10

Our goal is to compare the cardinality of the outputs from two phases of the Partitioning algorithm; at the Local level and the Global level, before and after application of our pre-processing. Firstly, data set is partitioned into fragments; secondly the Apriori algorithm (by Zhu T. [12]) is applied to find *Local frequent itemsets* (LP_1) for each fragment. Subsequently, union of these LP_i generates the *Global candidates*.

Resulting figures for each data set are represented in following template table 2. The 2nd, 3rd and 4th columns' names indicate three techniques for data preparation: *Sequent* fragments correspond to loading clusters with original data, *Skipping* fragments are constructed as described in section 5.1; and the *Clustering* fragments are the pre-processed data as presented by our clustering method described in section 5.2.

The data sets used are indicated on the top of each table segment. We present three different scenarios; each data set is partitioned into 1, 2 and 5 fragments. The *Sequent* column represents the numbers of the Local level (LP_1, LP_2, \dots, LP_n), the number of *Global candidates*. Note that this figure is presented by showing its two components; for example, **16 + (378)** indicates that there are **16** candidates to be checked and **378** common candidates don't need additional check.

Using the same convention, the *Skipping* and *Clustering* columns represent the figures for the *Skipping* technique and the *Clustering* pre-processing, respectively.

As can be seen from Table 2 and 3, there are big gains from the careful data pre-processing. Further, to discuss the impact of threshold level, let us denote the cardinality of checked Global candidate set as $|C_n^G|$, where n is the number of fragments. $|C_n^G|$ is reduced for all data sets for all support thresholds. For example, if T10I4D100k is partitioned into 2 fragments, $|C_2^G|$ decreases from **16** for *Sequent* to **3** for *Clustering* with the support threshold **0.01**. This reduction is also present when considering other real data sets that are partitioned into 2 fragments. Its value reduces from **1,820** to **348** with the threshold **0.005** for very large data set BMS-POS. Moreover, if data sets are partitioned into 5 fragments, this gap among 3 techniques is even greater. For example, if T10I4D100k is partitioned into 5 fragments, $|C_5^G|$ decreases from **48** for *Sequent* to **24** for *Clustering* with the threshold **0.01**, and **698** to **373** with

Table 2. The figures with a threshold 0.01

	Sequent	Skipping	Clustering
T1014D100K			
1-fragment: 385 Frequent Itemsets			
2 fragments			
LP1	385	387	385
LP2	387	386	386
C_2^G	16+ (378)	17 + (378)	3 + (384)
5 fragments			
LP1	392	386	387
LP2	381	388	387
LP3	393	388	384
LP4	386	387	388
LP5	390	391	388
C_5^G	48+ (366)	57 + (362)	24+ (375)
WebView-1			
1-fragment: 208 Frequent Itemsets			
LP1	227	241	210
LP2	229	201	213
C_2^G	152+ (152)	116 + (163)	17+ (203)
5 fragments			
LP1	284	250	226
LP2	197	230	221
LP3	241	254	213
LP4	255	242	207
LP5	266	205	205
C_5^G	425+ (92)	228 + (141)	74+ (181)
WebView-2			
1-fragment: 186 Frequent Itemsets			
LP1	279	156	192
LP2	221	236	179
C_2^G	292+(104)	120 + (136)	19+ (176)
5 fragments			
LP1	133	197	188
LP2	558	182	209
LP3	384	184	193
LP4	244	180	169
LP5	227	247	195
C_5^G	756+ (55)	157 + (135)	64+ (160)
BMS-POS			
1-fragment: 1,503 Frequent Itemsets			
LP1	1,400	1,353	1,512
LP2	1,662	1,680	1,498
C_2^G	390+ (1,336)	341+ (1,346)	60+ (1,475)
5 fragments			
LP1	1,996	1,719	1,150
LP2	1,334	1,146	1,471
LP3	744	1,639	1,864
LP4	1,348	1,810	1,822
LP5	2,885	1,377	1,364
C_5^G	2,263+ (689)	950+ (1,067)	894+ (1,121)

Table 3. The figures with threshold 0.005

	Sequent	Skipping	Clustering
T1014D100K			
1-fragment: 1,073 Frequent Itemsets			
2 fragments			
LP1	1,079	1,101	1,068
LP2	1,101	1,077	1,092
C_2^G	158 + (1,011)	148 + (1,015)	70 + (1,045)
5 fragments			
LP1	1,150	1,181	1,089
LP2	1,141	1,074	1,110
LP3	1,248	1,091	1,059
LP4	1,110	1,122	1,135
LP5	1,120	1,135	1,098
C_5^G	698 + (893)	578 + (889)	373 + (941)
WebView-1			
1-fragment: 633 Frequent Itemsets			
LP1	644	774	659
LP2	755	612	641
C_2^G	503 + (448)	416 + (485)	94 + (603)
5 fragments			
LP1	1,107	771	779
LP2	489	842	733
LP3	839	941	676
LP4	894	769	663
LP5	977	517	597
C_5^G	1,806 + (271)	1,069 + (374)	497 + (493)
WebView-2			
1-fragment: 996 Frequent Itemsets			
LP1	1,980	738	1,064
LP2	1,058	1,422	941
C_2^G	2,150 + (444)	808 + (676)	191 + (907)
5 fragments			
LP1	682	1,130	1,067
LP2	8,546	997	1,355
LP3	2,899	911	986
LP4	1,271	957	791
LP5	1,257	1,412	1,069
C_5^G	10,007+(230)	1,114 + (625)	751 + (723)
BMS-POS			
1-fragment: 6,017 Frequent Itemsets			
LP1	5,419	5,311	6,024
LP2	6,709	6,729	5,972
C_2^G	1,820+ (5,154)	1,468+ (5,286)	348+ (5,824)
5 fragments			
LP1	8,480	7,014	4,339
LP2	4,975	4,290	5,932
LP3	2,541	6,619	7,530
LP4	5,177	7,315	7,443
LP5	12,755	5,287	5,289
C_5^G	10,718+ (2,346)	4,353+ (3,956)	4,075+ (4,191)

the threshold **0.005**, respectively. Exceptional performance for WebView-2 data set with the threshold **0.005** the reduction is from **10,007** to only **751** when data set is partitioned into 5 fragments.

Hence naturally, another interesting and encouraging trend can be found in the growth of the number of common candidates between LP_i for fragmented data sets. For example, if data sets are partitioned into 5 fragments, this common number increases from **689** to **1,121** for BMS-POS with the threshold **0.01** as well as from **230** to **723** for WebView-2 with the threshold **0.005**.

In summary, the figures from 2 above tables show that the *Clustering* pre-processing technique can significantly improve the Partitioning approach. It is delivered in form of two strongly related benefits; reduction of the number of *Global candidates* requiring the final check and increase of the common candidates numbers that don't require any additional checks.

7 Conclusion

This paper considers a new approach for further performance improvements in frequent itemsets computation. Based on the original Partition algorithm, we show that the composition of fragments and the number of fragments generated, impact on the size of the data used by this algorithm.

We propose a pre-processing method (an incremental clustering algorithm), mainly to demonstrate that there is potential in the direction of performance improvement. Figures from the experiments show that this pre-processing offers good benefits already. The main question which still deserves consideration is related to the identification of methods that will deliver an even better partition for the original data sets.

Acknowledgment. We wish to thank the Data Mining group at ITEE School - The University of Queensland and the anonymous reviewers for suggestions.

References

- [1] Agrawal R., Srikant R.: *Fast algorithms for mining association rules*. Proc. 20th Int. Conf. Very Large Data Bases, Morgan Kaufmann, 1994 (487 - 499)
- [2] Savasere A., Omiecinski E., Navathe S.: *An efficient algorithms for mining association rules in large database*. Proc. 21th Int. Conf. Very Large Data Bases, Swizerland, 1995
- [3] Goethals B.: *Survey on frequent pattern mining*. University of Helsinki, 2002
- [4] Lin J.L., Dunham M.H.: *Mining association rules: Anti-skew algorithms*. Proc. 14th IEEE Int. Conf. on Data Engineering, Florida, 1998
- [5] Agrawal R., Imielinski T., Swami A.N.: *Mining association rules between sets of items in large database*. Proc. 1993 ACM SIGMOD Int. Conf. on Management of Data, 1993
- [6] Brin S., Motwani R., Ullman D.J., Tsur S.: *Dynamic Itemset Counting and implication rules for market basket data*. Proc. ACM SIGMOD 1997 Int. Conf. on Management of Data, 1997 (255 - 264)
- [7] Zaki M.J.: *Scalable algorithms for association mining*. IEEE Transactions on Knowledge and Data Engineering, 12(3): 372-390, 2000

- [8] Han J., Pei J., Yin Y., Mao R.: *Mining frequent patterns without candidate generation: A frequent-pattern tree approach*. Data Mining and Knowledge Discovery, Kluwer Academic Publishers, (8): 53-87, 2004
- [9] Mueller A.: *Fast sequential and parallel algorithm for association rules mining: A comparison*. Technical Report CS-TR-3515, University of Maryland, 1995
- [10] Son N. Nguyen: *Data partitioning approach into selected data mining problems*. PhD Confirmation report, The University of Queensland, Australia, 2005
- [11] Pudi V., Haritsa J.: *ARMOR: Association rule mining based on Oracle*. Workshop on Frequent Itemset Mining Implementations (FIMI'03 in conjunction with ICDM'03), 2003
- [12] Zhu T.: *The Apriori algorithm implementation*, <http://www.cs.ualberta.ca/~tszhu/>
- [13] Ron Kohavi, Carla Brodley, Brian Frasca, Lew Mason, and Zijian Zheng. *KDD-Cup 2000 organizers' report: Peeling the onion*. SIGKDD Explorations, 2(2):86-98, 2000

On-Line Adaptive Filtering of Web Pages

Richard Nock¹ and Babak Esfahani²

¹ GRIMAAG, Université Antilles-Guyane, Schoelcher, France
rnock@martinique.univ-ag.fr

² Dept of Systems and Computer Engineering,
Carleton University, Ottawa, Canada
babak@scce.carleton.ca

Abstract. We present a browser extension to dynamically learn to filter unwanted Uniform Resource Locators (such as advertisements or flashy images) based on minimal user feedback. Our extension builds upon one of the top ten of Mozilla Firefox plug-ins which filters URLs *without* learning capabilities. We apply a weighted majority-type learning algorithm working on regular expressions. Experimental results confirm that the accuracy of the predictions converges quickly to very high levels, with other key parameters: recall, specificity and precision.

1 Introduction

Many attempts have been made to make Web browsing more pleasant by allowing the user to remove big pictures and unwanted animations that interfere with reading. Some browsers such as Netscape or Mozilla allow the user to collapse such pictures or even create backlists of interfering domains that suppress them.

But the most sophisticated approaches so far have been proposed by the developers of AdBlock. AdBlock [1] is, according to Mozilla update data [6], the top ten of the most popular extensions to the Mozilla Firefox browser [5], with about 100000 downloads. To use AdBlock, the user has to come up with a collection of regular expressions that describe the URL patterns of images that they want to see filtered. As a result, whenever the browser is presented with a new URL, if the URL is matched by a regular expression, it is simply ignored, which does not necessarily speed up the web page, but at least makes page downloading faster.

```
/[a-z\d+%]{\w*\d+x\d}?(\d*(show)?(\w{3,}\%20|alligator|avs|barter|blog|box|central|d?html|i?frame|front|fuse|get|house|inline|instant|live|main|net|partner|primary|provider|rotated?|secure|side|smart|sponsor|story|text|view)?_?ads?(v?(bot|brite|broker|bureau|butler|center|click|client|creative|content|count|c|l|t)|data|engage|ler(tisw*|t(pro)?|ve|r?)|farm|force|frame|gif|group|id|head|id|img?ge?|info|js|juggler|legend|link|log|man(ager)?|max|mentor|meta\.com|net|optimi|sz|er|pic|popup|proof|q\.nextag|quest\.nl|redire?c?c?|remote|revolver|rotator|sale|sdk|sfac|solution|sonar|source|space|srv|stat.*\.asp|sys|track|trix|view|t?pe|zone)?)?(\d*(s|status)?\d*[W_]?(!*\w+\.edu|aware|adur1=|block|login|nl/|.*(&sb|(\.(\w|\r|\n))))
```

Fig. 1. Example of a long regular expression found on AdBlock's forum

However, as frequently discussed in the AdBlock developer forum, coming up with regular expressions is a difficult task, especially for the non-computer scientist. Writing and mastering them accurately requires extensive reading [2], and these published regular expressions can be especially hard to read and understand. Figure 1 presents

the example of a regular press. posted. AdBlock's forum. Most of the regular press. posted are smaller than this one, but some of them appear to be much more complicated to understand. Thus, the user faces the risk of obtaining unwanted browsing/blocking behaviors, sometimes without realizing them. The problem cannot be solved from a global standpoint, as it would be impossible to come up with a general set of filters that would satisfy every user. Finally, as the advertisement suppliers and browsing habits change, should the set of regular press. posts that are needed. The behavior of AdBlock is too static to be suited to these dynamic interactions, but making a dynamic interaction between the user and the filter is either thing but trivial. The conflict between the user has to be greater than its effective drawbacks, and the complexity of the algorithm is clear such a potential drawback.

To address this problem, we propose a fast machine learning approach that would create filters based on minimal interaction with the user. The user is not required to know how to create regular press. posts; all that is required is for the user to click on URLs (or images) that he/she wants to see blocked. Conversely, from time to time the user is needed to unblock URLs that should not have been blocked by the adaptive filter. Based on this simple feedback, our proposed method, an adaptation of the well-known Weighted Majority algorithm [4], builds a set of experts (simple regular press. posts or URLs) that determine whether a given URL should be blocked or not.

In the following Section, we review some works related to our topic. Then, the next Section is devoted to a formal presentation of the algorithm. After a Section presenting the browser extension, a Section presents and discusses experimental results. A last Section concludes and presents related issues in the topic.

2 Related Work

Our approach is inspired by the concept of Interface Agents [3]. An interface agent is a piece of software that assists a user of a complex system by observing his/her behavior and detecting repeating patterns that it could reproduce in order to automate tedious tasks. Typically such programs use some kind of machine learning algorithm to build the knowledge base. [3] designed interface agents that used k -nearest neighbor classification to share the filters with other agents. In a previous work, we used an adaptation of the Versatile Spaces algorithm to automate simple expert management tasks [7].

But the closest work is perhaps the use of Bayesian filtering for detecting email spam [8], which is now a standard feature in mainstream email programs. Bayesian methods for filtering emails have the advantage of being conceptually simple, and a great body of previous work has made them tailored to complex text classification tasks.

In our case, however, the setting makes them difficult to use as the best classification is suited for browsing. Classification is indeed made on a large scale. This is a crucial remark because the frequency of browsing through URLs is much higher than that of email receipt for the average user. This makes it necessary to have

an ultra-fast classification method with easy updates of the classifier, to filter the URLs as they come. In the case of email spam detection, it is a real necessity to have efficient feature selection algorithms to reduce the feature space to a small set prior to using Bayesian methods [8]. Making the addition of a hundred or more updates for URL filtering, such as the computation of the probability table for each feature, would rapidly slow down the browser and make its use unbearable. Furthermore, Bayesian methods rely on independence assumptions (at best partially relaxed) on the features to make the classification scalable [8]. This is clearly not a desirable assumption for URL classification, since it partially prevents the possibility of detecting a URL.

3 Theoretical Setting

Verifying formally, the algorithm can be reduced to the following finite steps: get a set of URLs, update a set of parameters, and update the classification of each element. A time during the algorithm, a prediction is possible for a set of URLs using a set of parameters over the current set of elements.

More formally, each browser address belongs to a set X , which contains a possible browser address. Each browser address is a URL (Uniform Resource Locator). From the user's standpoint, X can be partitioned into two subsets. The first one contains the URLs he would like to block, and refrain from adding. The other one contains the other URLs, which he wishes to leave unblocked. To each URL can thus be associated a status which we call β (block/unblock), and our objective is to predict the class of each URL as accurately as possible with respect to the user, given that a total of d different users may probably correspond to different partitions of X . Our algorithm builds therefore a decision function (classifier) from X to $\{-1, +1\}$, with $+1$ denoting the class of the URLs to be blocked (also called the "spam" class), and -1 the class of the URLs to leave unblocked (the "good" class).

We denote a couple (browser address, class) obtained from the user as a pair (x, y) . We let $(x_1, y_1), (x_2, y_2), \dots$ denote the stream of examples observed from the user, and (x_t, y_t) is thus the t^{th} example of the stream. We build a set of elements \mathbf{E} which is growing with time; to keep notations clear, we do not use the time subscript. \mathbf{E} : it should be clear from context which set of elements we use. Each element of \mathbf{E} is a couple (hypothesis, weight). A hypothesis is a function $h : X \rightarrow \{-1, 0, +1\}$ which is a prediction (this is the output 0). More precisely, each hypothesis' output is either $\{-1, 0\}$ or $\{0, +1\}$, which means that the corresponding element is authorized to say "I don't know", thus delegating the decision to the classifier associated to the other elements. The weight associated to hypothesis h is denoted $w_t(h) \in \mathbb{R}^+$. It is a function of time since it is updated each time an example is received. At the very beginning of the algorithm, prior to seeing the first example, we initialize the following set of parameters:

- $\beta \in (0, 1)$ is a personal static choice by the user,
- $\mathbf{E} \leftarrow \emptyset$ is the initial set of elements,
- $t \leftarrow 1$ is the time stamp abiding the examples received.

Algorithm 1 below displays a summary of what happens here. Example (x_t, y_t) is received.

Algorithm 1: Receive_New_Example $((x_t, y_t))$

Input: example (x_t, y_t)
 $\mathbf{N} \leftarrow \text{Create_Hypotheses}((x_t, y_t));$
 Update_Experts(\mathbf{N});
foreach $(h, w_t(h)) \in \mathbf{E}$ **do**
 $w_{t+1}(h) \leftarrow w_t(h) \times u(\beta, h, t);$
 $t \leftarrow t + 1;$

There are two possible choices for function $u(\beta, h, t)$:

$$u(\beta, h, t) = \frac{1 + y_t h(x_t)}{2\beta} + \frac{(1 - y_t h(x_t))\beta}{2}. \quad (1)$$

There are two procedures in Algorithm 1. **Create_Hypotheses**(.) takes an example as input, and outputs a set of hypotheses (regular expressions). Since the theoretical design of the algorithm does not depend on this procedure, we postpone the details and its implementation to the experimental section.

Update_Experts(.) takes as input a set of hypotheses, and creates a set of experts which is used to generate \mathbf{E} . In other words, it initializes the weights of the hypotheses. Details are given in Algorithm 2 (here, 0 denotes the function which is zero everywhere in \mathbb{R}).

Algorithm 2: Update_Experts(\mathbf{N})

Input: hypothesis set \mathbf{N}
foreach $h \in \mathbf{N}$ **do**
 $w_t(h) \leftarrow (u(\beta, 0, t))^{t-1};$
 $\mathbf{E} \leftarrow \mathbf{E} \cup \{(h, w_t(h))\};$

Weight initialization for the experts makes it possible to consider from the theoretical standpoint that each of them was created at the beginning of the algorithm, as if there were abstract universal input to \mathbf{E} . There remains to give the set \mathbf{E} is used to classify a user at $x \in X$. Just prior to receiving example $t + 1$, the decision made out of \mathbf{E} , $H_{\mathbf{E}, t}$, relies on a linear combination of: $\forall x \in X, H_{\mathbf{E}, t}(x) = \text{sign}(\sum_{(h, w_t(h)) \in \mathbf{E}} w_t(h) \times h(x))$.

4 Design of the Browser Extension

The Mozilla Firefox Web browser [5] is an open source product with an architecture specifically designed for a single 3rd party extension. This makes it

possibilities, especially the browser behavior, regarding the existing UI components, intercepting and reacting to browser events, and accessing relevant variables. Our findings regarding the test drivers are both implemented as such elements in JavaScript.

4.1 User Interface Elements

As a principle, a browser interface agent must remain as unobtrusive as possible, and therefore the user interface additions were kept to a minimum. We have provided the following items in the browser's content menu:

- the `cached_block_menu` which appears in the user right-clicks on a URL (or a image) that he/she wishes to block;
- the `other_cached_urls_block_menu` which is available should the user click on a URL that appears to be blocked by mistake. Selecting this item brings up the list of blocked items for the page, and the user can then choose which URL needs to be unblocked.

The `Block Menu` button is the one the user provides the positive response to the algorithm, while the `Unblock` button provides the negative ones. One could envisage that the `Unblock` items that were previously classified as such should also be fed to the algorithm (once the user has left the given page, thus confirming that the error correct itself unblocked) for eight refinement purposes, but we have decided against it, as we thought that if the user is the one triggering the response, he/she should have a better feeling that is happening behind the scenes. This remark also holds for the blocked items that were previously unblocked by the user. Finally, this allows the user to create and delete preferences, which themselves do not benefit the browser that much. Notice that updates of the preferences occur whenever receiving misclassified responses: false positives decrease the eight false positive errors (with regard to the `Block` class), while false negatives decrease the eight false negative errors (with regard to the `Unblock` class).

4.2 Implementation of `Create_Hypotheses(.)`

To generate the set of features N_1 Algorithm 1, we tokenize the response URLs using the character `/` as delimiter. The tokenized representations such as domain names, folders, but exclude file names. In that last case indeed, file names are frequently generated automatically for the URLs to block (or bad advertisement sites), and the resulting file names generally have little significance. Furthermore, this helps to keep the list simple and manageable. This is a simple choice of tokenization seems to be chosen by a significant proportion of users sharing their regular pressions. AdBlock's forums. Notice that `http` is a standard result key. The user may receive eight as the balance between the rate of false positives and the rate of false negatives achieved through earling, or, similarly, as a indication of the ratio between precision and recall.

The obtained keys are then compared with the corresponding set of features. Based on the corresponding mean that the keys obtained from positive (resp. negative) examples are compared to the positive (resp. negative) set of features. If a match is found, the new key is added to the corresponding set of features, and its weight is initialized using Algorithm 2. More keys could be used to be generated. For instance, we could also use the full URL itself as a feature. Also, the character could be used as a delimiter, the periods of the parts of a domain name that are kept to its classification (e.g. homepage.com). Finally, we could create features that capture the importance of the most significant keys appearing in a URL. The factors to consider here are to avoid the proliferation of features.

5 Experimental Results

In our experiments, we have fixed $\beta = 1/\sqrt{e} \approx 0.61$ in update rule (1). In order to obtain results that are independent from a particular browsing habit, we needed to provide a test setting that could be used seamlessly by a kind user. Therefore, in addition to providing the standard evaluation described in the previous section, we embedded our algorithm inside the AdBlock extension.

The AdBlock user is asked to set up filters as usual in the form of regular expressions, creating as a result a race for the embedded ear. The AdBlock filters override the ear's classification in order to remain transparent to the end user. This means that to the user, the extension is behaving differently than the regular AdBlock. However, an earer misclassification (false positives and false negatives) are fed back as such to the algorithm, leading to the new feature creation and weight adjustments described above.

At each step consisting of k browser actions (e.g. visited image URLs), we freeze a copy of the earer's knowledge base up to that point. While the earer keeps evaluating and accepting feedback from the race, the copy is used to evaluate the earer's accuracy of the accumulated knowledge so far by applying a confusion matrix based on its predictions to the corresponding examples. After n such steps, and for a total of $n \times k$ browser actions, the user is notified that the testing is finished, and the logs are collected. We can therefore compare the earer at each step and observe the evolution of its ability to classify the upcoming browser actions. However, as we get close to the final steps of each test, the number of browser actions available to the more recent earers decreases, and the statistical confidence in the more recent results decreases as well. To reduce this phenomenon, we allow some more browser actions to be collected after the last step.

5.1 AdBlocking on a Single Commercial Website

Our first set of tests were designed to see whether our algorithm is able to correctly predict which URLs to block (e.g. business-related images) on a page, and if so, after how many visits. We used a commercial

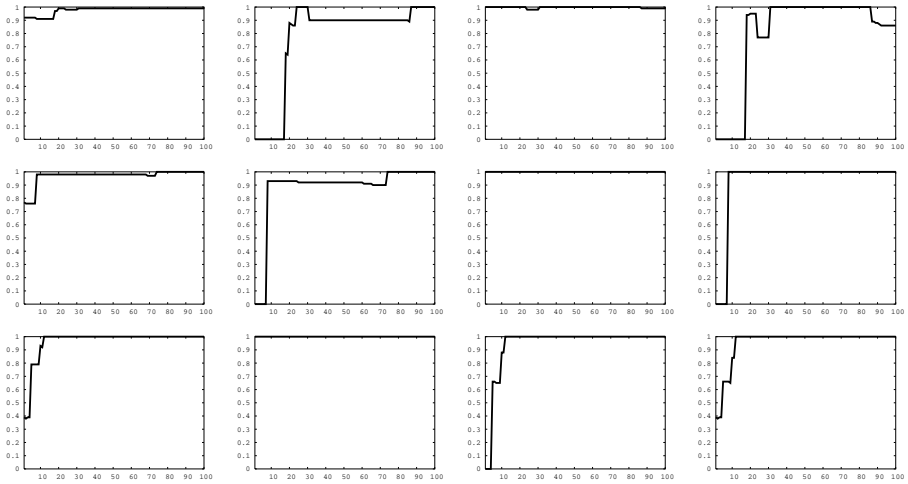


Fig. 2. From left to right: accuracy, recall, specificity and precision. From up to bottom: websites of CNN, Fox News and MSN (x -axis=step number, see text for details).

set of AdBlock regular expressions, as compared to the used AdBlock discussion forums, as race. On three popular and large commercial websites, we have run AdBlockLearn with $k = 1$ and $n = 100$. The total number of users who were usually reached is quick. Figure 2 plots the evolution of four parameters through out early stage. If we denote by TP the number of true positives, TN the number of true negatives, FP the number of false positives and FN the number of false negatives, then the accuracy is $TN/(TN + FP)$, the recall is $TP/(TP + FN)$, and the precision is $TP/(TP + FP)$. As can be seen, the algorithm converges quickly to a good prediction in terms of four parameters. This is good given that commercial websites use dynamic advertising advertisements using cookies, and as a result hitting relevant users by a direct set of images and URLs. However it is important to point out that the total similar results in a k -test setting, the miscellaneous that were detected by the race would have to correspond to as many direct feedbacks by the user. In practice, in the absence of similar interactions, the four parameters can be sub-optimal, but it is definitely acceptable.

5.2 AdBlocking While Surfing to Different Websites

The set of tests measures robustness to overfitting. How does the early detection of edge-tracker on other websites, are the rules learned so far useful to the websites, and how much more early going is left to do? Our intuition is that the amount of miscellaneous should decrease over time, as usually the providers of false advertisement are the same in many different commercial sites. We asked the users to simplify their usual browsing habits, and e

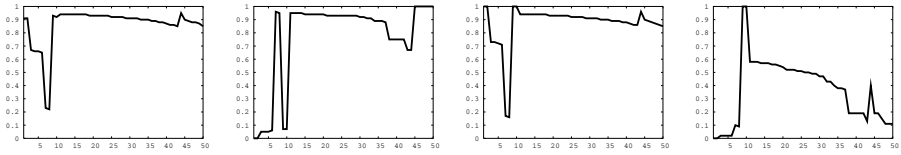


Fig. 3. Evolution of the four parameters during a typical browsing session (conventions follow figure 2, see text for details)

set $k = 10$ (to absorb some of the variability) and $n = 50$. To make the task harder, we requested that each chosen web site had to be visited \dots . The results are charted in figure 3. It is quite remarkable that the accuracy, the specificity, the recall and the precision remained at very high values after a short period, given the tough parameter setting. However, the fact that the precision decreases tends to indicate that there is a significant increase in FP (to make the precision decrease) and TN (to make the specificity remain at high values). This may dispel the fact that the number of negative examples does not increase, but the experts might be too simple to fit the growing amount of information, to discriminate among the examples that come from various sources. In case false positives are deemed unacceptable by the user, i.e. the user does not want to have to manually check erroneous blocked URLs, it is possible that after the utilization of the weighted majority algorithm more negative examples, as a side effect [8]. The trade-off would be a drop in accuracy and a slight increase in the rate of false negatives.

6 Conclusion and Future Work

In this paper, we have experimentally demonstrated the efficiency of a carefully adapted weighted majority. Compared to usual weighted majority, our setting makes use of the fact that experts may abstain instead of making a guess. This raises an important technical issue, as the efficiency of weighted majority is usually measured with respect to its number of \dots [4]. In our setting, we would certainly appreciate this quality to be as small as possible, \dots we would appreciate the number of abstentions to be small. Since mistake bounds do not take into account the number of abstentions, this raises both the problem of finding accurate qualities to minimize, and relevant bounds that our adapted weighted majority satisfies.

Acknowledgments and Code Availability

We would like to thank Rue, chief developer of AdBlock, for his time and enthusiasm, and for recommending us to his team. R. Nock would like to thank Ottavia U. University and StatMat for financial support, during which part of this work was achieved. Both the standard IEEE style

(AdBlockLearner) and the test driver (AdBlockLearnerTest) are available at <http://adblocker.mzdev.org>, including source code and documentation. They are compatible with most versions of Mozilla Firefox.

References

1. AdBlock, 2005. <http://adblock.mozdev.org>.
2. Jeffrey E. F. Friedl. *Mastering Regular Expressions*. O' Reilly, 1997.
3. Y. Lashkari, M. Metral, and P. Maes. Collaborative interface agents. In *Proc. of AAAI-94*, pages 444–449, 1994.
4. N. Littlestone and M. Warmuth. The weighted majority algorithm. *Information and Computation*, pages 212–261, 1994.
5. Mozilla Firefox, 2005. <http://mozilla.org/products/firefox>.
6. Mozilla Firefox Extensions, 2005. <http://update.mozilla.org/extensions/>.
7. R. Nock and B. Esfandiari. Oracles and assistants : machine learning applied to network supervision. In *Canadian Artificial Intelligence Conference*, number 1418 in Lecture Notes in Computer Science, pages 86–98. Springer-Verlag, 1998.
8. M. Sahami, S. Dumais, D. Heckerman, and E. Horvitz. A bayesian approach to filtering junk email. In *AAAI Workshop on Learning for Text Categorization*. AAAI Press, 1998.

Table 1. A Boolean context \mathbf{r}

	g_1	g_2	g_3	g_4	g_5
t_1	1	0	1	1	0
t_2	0	1	0	0	1
t_3	1	0	1	1	0
t_4	0	0	1	1	0
t_5	1	1	0	0	1
t_6	0	1	0	0	1
t_7	0	0	0	0	1

... e a 1 g bi-c ... F... a c... e a a ... [6] ... h be a ... 1 ... I f ...
 ... a , a f... a c... e i a bi-e (T, G) he e he e f... b ec... T a d he
 ... e f... e ie G f... a a a (c... bi a... ia) , e c a g e f... e a e ,
 ... e.g., $(\{t_1, t_3\}, \{g_1, g_3, g_4\})$ i r. U f... a e , e g e e a g e h g e c e c i...
 ... f... a c... e ... h c h a e d i c... i e , e b e d - e e . I... e a e ,
 $\{\{t_1, t_3, t_4\}, \{g_1, g_3, g_4\}\}$ i... a f... a c... e $((t_4, g_1) \notin \mathbf{r})$ b... c a b e b i
 f... $\{\{t_1, t_3\}, \{g_1, g_3, g_4\}\}$ a d $\{\{t_1, t_3, t_4\}, \{g_3, g_4\}\}$ h... c h a e i i a e... g h
 f... a c... e . I... d e h e i... i l... f... a... a c h.

The c... i b i... f h... a e i... f d. F l... e... e a e b i-c... e i g
 f a e... h... h e a b e... c... e b i- a... i l... b g... i g... c a... a e...
 h... c h a... e... c a... g a... c i a... b e e... b e c... a d... e... i e , i.e.,
 b i- e... h... c h a i f... e... e- d e d c... a l... V a... c a... a e... a e
 c a d i d a e f... c h a... c e... e.g., f e... e... f... e... i e a... c i a e d... h e
 ... i g e... f... b e c... f... a c... e... , e c. S e c... d... e... d... e i... a c e f
 h... f a e... , h e C D K - M E A N S a g... i h... , h... c h b i d... i... a e... i e d
 ... a... i l... b e c... a d... e... i e . M... e... e c i e... , e a... a K - M E A N S - i e
 a g... i h... a c... e c i... f b i- e... (f... a c... e... i... , e... e... i e...). A a
 ... e... , b e c... a d... e... i e a e i... i... c a... a... c i a e d... c... e... , d e d i g
 ... h e... e i g h... i... h e... a... c... e d c... e... i d . O... e... e... i e a... a i d a i...
 c... e... h e a d d e - a... e f C D K - M E A N S... h e (b i) c... e i g a g... i h... .
 I... S e c... 2... , e... e... c... e i g f a e... , a d... e... e... e a e d
 S e c... 3 d i c... e... , e... e... i e a... a i d a i... . e h d... g a d i c... -
 a... a... e... e... i e a... e... a... i... b e c h a... d a a e... . A c... a... i...
 b e e... C D K - M E A N S... , b i- c... e i g a g... i h... . (C O C L U S T E R [4] a d B i-
 C L U S T [3]), a d... c a... i c a... c... e i g a g... i h... . (K - M E A N S a d E M [1]) i
 g i e . S c a a b i... i... e a e d i c... e d a d S e c... 4 c... c... d e .

2 Clustering Model

A... e a e... f... b e c... $\mathcal{O} = \{t_1, \dots, t_m\}$ a d a e... f... B... e a... e... i e $\mathcal{P} =$
 $\{g_1, \dots, g_n\}$. The B... e a... c... e... b e... i e d i r $\mathbf{r} \subseteq \mathcal{O} \times \mathcal{P}$, h e e $r_{ij} = 1$ f
 ... e... g_j i... a i e d b... b e c... t_i . We d e... e h e b i- c... e i g a... a f... :
 ... e a... c... e... e a... a... i l... f K c... e... f... b e c... (a... $\{C_1^o \dots C_K^o\}$) a d a
 ... a... i l... f K c... e... f... e... i e (a... $\{C_1^p \dots C_K^p\}$) i h a... a... i g b e... e...
 b... h... a... i l... c h h a e a c h c... e... f... b e c... i c h a c a e i e d b a c... e... f
 ... e... i e . O... i d e a i... h a b i- a... i l... c a... b e c... e d f... b i- e... a d i... i

be 1. a ia ed a e . . . f . . . a c . ce . . . F . . . a , a bi-e 1 a e e e . $b_j = (T_j, G_j)$ ($T_j \subseteq \mathcal{O}, G_j \subseteq \mathcal{P}$) a d e a . . . e ha a c . ec 1 . . . f a . . . 1, 11 e e 1 g bi-e . de . ed \mathcal{B} ha bee e . ac ed f . . . r bef . e ha d. Le . . . de c, i be b_j b . he B . ea . ec . . . $\langle \mathbf{t}_j \rangle, \langle \mathbf{g}_j \rangle = \langle t_{j1}, \dots, t_{jm} \rangle, \langle g_{j1}, \dots, g_{jn} \rangle$ he e $t_{jk} = 1$ if $t_k \in T_j$ (0 . he . 1 e) a d $g_{jk} = 1$ if $g_k \in G_j$ (0 . he . 1 e). We a e . . . 1 g f . . . K c . e . . . f bi-e . $\{C_1, \dots, C_K\}$ ($C_i \subseteq \mathcal{B}$). Le . . . de . e he ce . . . id f a c . e . . . f bi-e . C_i a $\mu_i = \langle \tau_i \rangle, \langle \gamma_i \rangle = \langle \tau_{i1}, \dots, \tau_{im} \rangle, \langle \gamma_{i1}, \dots, \gamma_{in} \rangle$ he e τ a d γ a e he . . . a ce . . . id c . . . e . . . :

$$\tau_{ik} = \frac{1}{|C_i|} \sum_{b_j \in C_i} t_{jk}, \quad \gamma_{ik} = \frac{1}{|C_i|} \sum_{b_j \in C_i} g_{jk}$$

We . . . de . e . . . di a ce be ee a bi-e a d a ce . . . id:

$$d(b_j, \mu_i) = \frac{1}{2} \left(\frac{|\mathbf{t}_j \cup \boldsymbol{\tau}_i| - |\mathbf{t}_j \cap \boldsymbol{\tau}_i|}{|\mathbf{t}_j \cup \boldsymbol{\tau}_i|} + \frac{|\mathbf{g}_j \cup \boldsymbol{\gamma}_i| - |\mathbf{g}_j \cap \boldsymbol{\gamma}_i|}{|\mathbf{g}_j \cup \boldsymbol{\gamma}_i|} \right)$$

I 1 . he . ea . f he eigh ed . . . e, i ca di e e ce . f he e c . . . e . . . We a . . . e $|\mathbf{t}_j \cap \boldsymbol{\tau}_i| = \sum_{k=1}^m a_k \frac{t_{jk} + \tau_{ik}}{2}$ a d $|\mathbf{t}_j \cup \boldsymbol{\tau}_i| = \sum_{k=1}^m \frac{t_{jk} + \tau_{ik}}{2}$ he e $a_k = 1$ if $t_{jk} \cdot \tau_{ik} \neq 0, 0$. he . 1 e. I 1 1 e , he 1 e, ec 1 . 1 e a . . he . ea be ee he . . be . f c be c . a d he . . . f he c . . . id eigh . . The . 1 . 1 . he . ea be ee he . . be . f b ec . a d he . . . f he c . . . id eigh . . The e . ea . e a e de . ed . 1 1 a . . . e . e .

Ob ec . t_j (e . . . e . ie g_j) a e a i g ed . . . e f he K c . e . (de . . ed i) f . . hch τ_{ij} (e . . . γ_{ij}) 1 . . a 1 . . . We ca e a be ha a . . be . . f b ec . a d / . . . e . ie be . g . . . e ha . . . ec . e b c . . . i g he . 1 e f he . e a i g a . f each c . e . Tha de . 1 1 . f c . e . . e be . hi de e . 1 ed b . he a e . f $\boldsymbol{\tau}_i$ a d $\boldsymbol{\gamma}_i$, e . . . eed . ada he c . e a i g . e . e . F . . hi . . . e, e . . 1 . d ce a a e e . δ_o a d δ_p 1 $[0, 1]$. . a if he . e be . hi f each e e . . a c . e . We a ha a . b ec t_j be . g . . a c . e C_i^o if $\tau_{ij} \geq (1 - \delta_o) \cdot \max_i(\tau_{ij})$. A a g . . , a . . . e . g_j be . g . . a c . e C_i^p if $\gamma_{ij} \geq (1 - \delta_p) \cdot \max_i(\gamma_{ij})$. Ob 1 . . he . . be . f . e a i g . b ec . (e . . . e . ie) de e d . . he di . b 1 . . f he a e . f $\boldsymbol{\tau}_i$ (e . . . $\boldsymbol{\gamma}_i$). N . ice ha if . e a i g 1 a . ed, $\delta = 0$ de . . 1 . . . ha each b ec . . . e . 1 a i g ed . . a i g e c . e . The ch ice . f a . e e a . a e f . δ 1 ce a . a i ca 1 . - de e de . . Whe a bi-c . e 1 g . . c . e h d 1 . he da a , 1 e . a e . f δ a e . . e . gh . . . ide . e e a . . e a i g . O a . he ha d , 1 . 1 . c . e . , e e . 1 e a e . f δ ca g i e . 1 e . . i g 1 ca . . e a i g . e .

We ca ide de a . ab . . he . died 1 . a ce . f hi f a e . . : a bi-c . e 1 g ba ed . . f . . a c . ce . . Ma . e ce . a g , i h . ha e bee . de e . ed ha ca e . ac c . . e e c . ec 1 . . . f f . . a c . ce . . . de . c . . . a . . . We . e D-MINER [7].

O . 1 . a ce CDK-MEANS 1 . e e . ed 1 . Tab e 2. I c . . . e a bi- a . 1 1 . . f a da a e r g i e . a c . ec 1 . . f bi-e . \mathcal{B} e . ac ed f . . . r bef . e ha d (e.g., f . . a c . ce . .), he de 1 ed . . be . f c . e . K , he h e h d a e f . δ_o

Table 2. CDK-MEANS pseudo-code

CDK-MEANS (\mathbf{r} is a Boolean context, \mathcal{B} is a collection of bi-sets in \mathbf{r} , K is the number of clusters, MI is the maximal iteration number, δ_o and δ_p are thresholds values for controlling overlapping)

1. Let $\mu_1 \dots \mu_K$ be the initial cluster centroids. $k := 0$.
2. Repeat
 - (a) For each bi-set $c \in \mathcal{B}$, assign it to cluster C s.t. $d(c, \mu_i)$ is minimal.
 - (b) For each cluster C_i , compute τ_i and γ_i .
 - (c) $k := k + 1$.
3. Until centroids are unchanged or $k = MI$.
4. If overlap is allowed, for each $t_j \in \mathcal{O}$ (resp. $g_j \in \mathcal{P}$), assign it to each cluster C_i^o (resp. C_i^p) s.t. $\tau_{ij} \geq (1 - \delta_o) \cdot \max_i(\tau_{ij})$ (resp. $\gamma_{ij} \geq (1 - \delta_p) \cdot \max_i(\gamma_{ij})$).
5. Else, for each $t_j \in \mathcal{O}$ (resp. $g_i \in \mathcal{P}$), assign it to the first cluster C_i^o (resp. C_i^p) s.t. τ_{ij} (resp. γ_{ij}) is max.
6. Return $\{C_1^o \dots C_K^o\}$ and $\{C_1^p \dots C_K^p\}$

and δ_p , and a τ and a γ are defined by $\tau_{ij} = \max_i(\tau_{ij})$ and $\gamma_{ij} = \max_i(\gamma_{ij})$. Otherwise, CDK-MEANS provides the bi-assignment given in Section 1. The complexity is linear in K and MI . See also the related work in Section 3.

Related work. [3] and [4] bi-cluster algorithms deal with overlapping clusters. The algorithm proposed in this paper is based on the G-d algorithm [8]. The G-d algorithm is a greedy algorithm that iteratively assigns each bi-set to the cluster that minimizes the distance to the centroid. The algorithm is based on the G-d algorithm [8]. The algorithm is based on the G-d algorithm [8]. The algorithm is based on the G-d algorithm [8].

3 Experimental Validation

The experimental evaluation is based on the G-d algorithm [8]. The algorithm is based on the G-d algorithm [8]. The algorithm is based on the G-d algorithm [8]. The algorithm is based on the G-d algorithm [8].

The algorithm is based on the G-d algorithm [8]. The algorithm is based on the G-d algorithm [8]. The algorithm is based on the G-d algorithm [8]. The algorithm is based on the G-d algorithm [8].

$$\tau_Q = \frac{\sum_i \sum_j \frac{p_{ij}^2}{p_{i.}} - \sum_j p_{.j}^2}{1 - \sum_j p_{.j}^2}$$

We evaluated the efficiency of the eight algorithms on data sets available from the UCI ML Repository¹ and from the JSE Data Archive². All the experiments have been performed on a PC with 1 Gb RAM and a 3.0 GHz P4 processor. First, we considered the categorical data sets available in the D-MINER [7]. Moreover, the categorical data sets have been used from mushroom and credit-a (11 attributes and 6 clusters (13, 15) and (6, 15)) available from the eec.ec1.org website before using CDK-MEANS.

Table 3. Goodman-Kruskal’s coefficient values for different bi-clustering algorithms (MR-2 and MR-5 refer to mushroom with 2 and 5 clusters)

Dataset	Dim.	BI-CLUST	COCLUSTER		CDK-MEANS	
		Max	Max	Mean	Max	Mean
voting	435×48	0.320	0.320	0.315±0.002	0.311	0.311±0.000
titanic	2201×8	0.332	0.321	0.226±0.076	0.314	0.160±0.109
iris-2	150×8	0.543	0.543	0.357±0.195	0.543	0.474±0.056
iris-3	150×8	0.544	0.390	0.379±0.045	0.523	0.329±0.080
zoo-2	101×16	0.191	0.186	0.157±0.034	0.192	0.165±0.020
zoo-7	101×16	-	0.080	0.065±0.009	0.083	0.049±0.015
breast-w	699×18	0.507	0.507	0.474±0.121	0.498	0.498±0.000
credit-3	690×52	0.104	0.014	0.003±0.003	0.110	0.091±0.015
credit-2	690×52	-	0.012	0.006±0.004	0.096	0.055±0.011
mr-2	8124×126	-	0.198	0.158±0.026	0.176	0.157±0.017
mr-5	8124×126	0.187	0.119	0.097±0.009	0.116	0.112±0.004
ads	3279×1555	-	0.006	0.003±0.001	0.538	0.137±0.109

We compared CDK-MEANS bi-clustering with the hybrid bi-clustering COCLUSTER [4], and BI-CLUST [3]. All the 11 data sets of the eight algorithms are divided into 100 clusters each data set and executed the algorithm which yielded the best Goodman-Kruskal’s coefficient. The best performed coefficients of each data set have been used to compare the categorical data sets. The coefficient of BI-CLUST which is a categorical data set is the best performed BI-CLUST 1 attribute in WEKA³ and the best performed categorical internet-ads (there are 1500 clusters). We compared the results in Table 3. We considered the τ_Q coefficient. The coefficient of τ_Q coefficient are a significant indicator. Notice that the CDK-MEANS has the best performance. The Goodman-Kruskal’s coefficient of the significant data sets of the eight algorithms. Overall, the hybrid internet-ads, the coefficient of the CDK-MEANS is comparable with the COCLUSTER. This demonstrates the high efficiency of the data sets which have been used by the eight algorithms. All the average behavior of the eight COCLUSTER. The average performance of the eight algorithms are as follows:

¹ <http://www.ics.uci.edu/~mllearn/MLRepository.html>
² http://www.amstat.org/publications/jse/jse_data_archive.html
³ <http://www.cs.waikato.ac.nz/ml/weka/>

Table 4. Jaccard coefficient values w.r.t. class variable for different algorithms

Dataset	BI-CLUST	COCLUSTER	K-MEANS	EM	CDK-MEANS
voting	0.6473	0.6473	0.6027	0.6459	0.6737
titanic	0.4281	0.4651	0.3697	0.3697	0.4745
iris-2	0.4992	0.4992	0.5117	0.4992	0.4992
iris-3	0.4932	0.5240	0.5394	0.5394	0.5144
zoo-2	0.5141	0.5630	0.5027	0.5179	0.5141
zoo-7	-	0.1647	0.1843	0.2325	0.2212
breast-w	0.8246	0.8287	0.7777	0.8328	0.7666
credit-3	0.4233	0.3869	0.3765	0.3405	0.4452
credit-2	-	0.4360	0.4698	0.4442	0.4915
mr-2	-	0.6819	0.3496	0.6976	0.6356
mr-5	0.5068	0.3450	0.3192	0.3364	0.3375
ads	-	0.4317	-	-	0.8019

the dataset of 10000 records. Notice that for voting-records and breast-w, CDK-MEANS has a Jaccard coefficient higher than BI-CLUST.

CDK-MEANS generally leads to the best performance for the datasets where the algorithm becomes more difficult to apply. For example, in the case of the titanic dataset, CDK-MEANS has a Jaccard coefficient of 0.4745, which is higher than the other algorithms. In the case of the iris-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.4992, which is higher than the other algorithms. In the case of the iris-3 dataset, CDK-MEANS has a Jaccard coefficient of 0.5394, which is higher than the other algorithms. In the case of the zoo-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.5179, which is higher than the other algorithms. In the case of the zoo-7 dataset, CDK-MEANS has a Jaccard coefficient of 0.2325, which is higher than the other algorithms. In the case of the breast-w dataset, CDK-MEANS has a Jaccard coefficient of 0.8328, which is higher than the other algorithms. In the case of the credit-3 dataset, CDK-MEANS has a Jaccard coefficient of 0.4452, which is higher than the other algorithms. In the case of the credit-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.4915, which is higher than the other algorithms. In the case of the mr-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.6976, which is higher than the other algorithms. In the case of the mr-5 dataset, CDK-MEANS has a Jaccard coefficient of 0.3375, which is higher than the other algorithms. In the case of the ads dataset, CDK-MEANS has a Jaccard coefficient of 0.8019, which is higher than the other algorithms.

We also tested the Jaccard index for the other algorithms. For example, for the titanic dataset, WEKA has a Jaccard coefficient of 0.3697, which is lower than the other algorithms. For the iris-2 dataset, WEKA has a Jaccard coefficient of 0.4992, which is equal to the other algorithms. For the iris-3 dataset, WEKA has a Jaccard coefficient of 0.5394, which is equal to the other algorithms. For the zoo-2 dataset, WEKA has a Jaccard coefficient of 0.5179, which is equal to the other algorithms. For the zoo-7 dataset, WEKA has a Jaccard coefficient of 0.1843, which is lower than the other algorithms. For the breast-w dataset, WEKA has a Jaccard coefficient of 0.7777, which is lower than the other algorithms. For the credit-3 dataset, WEKA has a Jaccard coefficient of 0.3765, which is lower than the other algorithms. For the credit-2 dataset, WEKA has a Jaccard coefficient of 0.4698, which is lower than the other algorithms. For the mr-2 dataset, WEKA has a Jaccard coefficient of 0.3496, which is lower than the other algorithms. For the mr-5 dataset, WEKA has a Jaccard coefficient of 0.3192, which is lower than the other algorithms. For the ads dataset, WEKA has a Jaccard coefficient of 0.4317, which is lower than the other algorithms.

First, we have seen that the Jaccard index is a good measure of the similarity between two datasets. In this paper, we have used the Jaccard index to compare the performance of different algorithms. We have seen that CDK-MEANS generally leads to the best performance for the datasets where the algorithm becomes more difficult to apply. For example, in the case of the titanic dataset, CDK-MEANS has a Jaccard coefficient of 0.4745, which is higher than the other algorithms. In the case of the iris-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.4992, which is higher than the other algorithms. In the case of the iris-3 dataset, CDK-MEANS has a Jaccard coefficient of 0.5394, which is higher than the other algorithms. In the case of the zoo-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.5179, which is higher than the other algorithms. In the case of the zoo-7 dataset, CDK-MEANS has a Jaccard coefficient of 0.2325, which is higher than the other algorithms. In the case of the breast-w dataset, CDK-MEANS has a Jaccard coefficient of 0.8328, which is higher than the other algorithms. In the case of the credit-3 dataset, CDK-MEANS has a Jaccard coefficient of 0.4452, which is higher than the other algorithms. In the case of the credit-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.4915, which is higher than the other algorithms. In the case of the mr-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.6976, which is higher than the other algorithms. In the case of the mr-5 dataset, CDK-MEANS has a Jaccard coefficient of 0.3375, which is higher than the other algorithms. In the case of the ads dataset, CDK-MEANS has a Jaccard coefficient of 0.8019, which is higher than the other algorithms.

Scalability Issues. Computing the Jaccard index for large datasets can be a challenging task. In this paper, we have used the Jaccard index to compare the performance of different algorithms. We have seen that CDK-MEANS generally leads to the best performance for the datasets where the algorithm becomes more difficult to apply. For example, in the case of the titanic dataset, CDK-MEANS has a Jaccard coefficient of 0.4745, which is higher than the other algorithms. In the case of the iris-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.4992, which is higher than the other algorithms. In the case of the iris-3 dataset, CDK-MEANS has a Jaccard coefficient of 0.5394, which is higher than the other algorithms. In the case of the zoo-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.5179, which is higher than the other algorithms. In the case of the zoo-7 dataset, CDK-MEANS has a Jaccard coefficient of 0.2325, which is higher than the other algorithms. In the case of the breast-w dataset, CDK-MEANS has a Jaccard coefficient of 0.8328, which is higher than the other algorithms. In the case of the credit-3 dataset, CDK-MEANS has a Jaccard coefficient of 0.4452, which is higher than the other algorithms. In the case of the credit-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.4915, which is higher than the other algorithms. In the case of the mr-2 dataset, CDK-MEANS has a Jaccard coefficient of 0.6976, which is higher than the other algorithms. In the case of the mr-5 dataset, CDK-MEANS has a Jaccard coefficient of 0.3375, which is higher than the other algorithms. In the case of the ads dataset, CDK-MEANS has a Jaccard coefficient of 0.8019, which is higher than the other algorithms.

⁴ Clearly, it does not lead to the highest Jaccard's index.

4 Conclusion and Future Work

We have introduced a new biclustering algorithm which is particularly suited to the data sets considered in this paper. The success of CDK-MEANS is due to the fact that it is able to find a good biclustering for a wide range of data sets. Our experimental results have shown that the added value of CDK-MEANS over the (bi-)clustering algorithms is that it is able to find a good biclustering for a wide range of data sets. The success of CDK-MEANS is due to the fact that it is able to find a good biclustering for a wide range of data sets. Our experimental results have shown that the added value of CDK-MEANS over the (bi-)clustering algorithms is that it is able to find a good biclustering for a wide range of data sets.

Acknowledgements. The authors would like to thank Ligia Madeira for her technical assistance. This research is partially funded by CNRS (ACI MD 46 B1 g).

References

- Jain, A., Dubes, R.: Algorithms for clustering data. Prentice Hall, Englewood cliffs, New Jersey (1988)
- Fisher, D.H.: Knowledge acquisition via incremental conceptual clustering. *Machine Learning* **2** (1987) 139–172
- Robardet, C., Feschet, F.: Efficient local search in conceptual clustering. In: Proceedings DS'01. Number 2226 in LNCS, Springer-Verlag (2001) 323–335
- Dhillon, I.S., Mallela, S., Modha, D.S.: Information-theoretic co-clustering. In: Proceedings ACM SIGKDD 2003, Washington, USA, ACM Press (2003) 89–98
- Madeira, S.C., Oliveira, A.L.: Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM Trans. Comput. Biol. Bioinf.* **1** (2004) 24–45
- Wille, R.: Restructuring lattice theory: an approach based on hierarchies of concepts. In Rival, I., ed.: *Ordered sets*. Reidel (1982) 445–470
- Besson, J., Robardet, C., Boulicaut, J.F., Rome, S.: Constraint-based concept mining and its application to microarray data analysis. *Intelligent Data Analysis* **9**(1) (2005) 59–82
- Goodman, L.A., Kruskal, W.H.: Measures of association for cross classification. *Journal of the American Statistical Association* **49** (1954) 732–764
- Pensa, R.G., Robardet, C., Boulicaut, J.F.: Using locally relevant bi-sets for categorical data conceptual clustering. Research report, LIRIS CNRS UMR 5205 - INSA Lyon, Villeurbanne, France (2005) Submitted to a journal.

Privacy-Preserving Collaborative Filtering on Vertically Partitioned Data*

Husein Polat and Weiqiang Du

Department of Electrical Engineering and Computer Science,
Syracuse University, CST 3-114, Syracuse, NY 13244-1240, USA
{hpolat, wedu}@ecs.syr.edu

Abstract. Collaborative filtering (CF) systems are widely used by E-commerce sites to provide predictions using existing databases comprised of ratings recorded from groups of people evaluating various items, sometimes, however, such systems' ratings are split among different parties. To provide better filtering services, such parties may wish to share their data. However, due to privacy concerns, data owners do not want to disclose data. This paper presents a privacy-preserving protocol for CF grounded on vertically partitioned data. We conducted various experiments to evaluate the overall performance of our scheme.

1 Introduction

Collaborative filtering (CF) is a recent technique that helps users cope with information overload using other users' preferences. It is widely used by E-commerce, direct recommendation systems, and search engines [1,2]. The goal is to predict how a user (e.g., u) will like a item that he/she did not buy before based on other users' preferences [4].

Data collected for CF purposes might be vertically partitioned between different parties where the parties hold disjoint sets of items' ratings collected from the same users. An individual's preferences for products might be split among different E-commerce companies such as Amazon.com and Microsoft. Online retailers can produce better referrals if they share information about their customers with their partners. Joint data is beneficial for E-commerce sites because customers prefer returning to sites with better referrals and the search for merchandise to purchase. Shared information is a significant benefit customers bring to get more accurate and reliable recommendations. Combining vertically partitioned data (VPD) is helpful because CF systems have limited rated items. Therefore, more reliable matchings and provide more accurate referrals, the cooperation between users should be encouraged; this might be achieved by integrating VPD. However, due to privacy concerns, data owners do not want to totally abandon and disclose their data to each other.

* This work was supported by Grants ISS-0219560 and ISS-0312366 from the United States National Science Foundation.

VPD-based CF is essentially achievable if privacy measures are introduced to data users. We studied the privacy-preserving collaborative filtering (PPCF) on VPD problem:

$$(A \rightarrow B)$$

We propose a practical achievable PPCF on VPD. Privacy, accuracy, and efficiency are conflicting goals. Therefore, the proposed practical should achieve a good balance between them. Our scheme consists of a $(n-1)$ and $(n-1)$ encrypted computation steps. We conduct some computation steps to achieve data exchange between parties with privacy. During the $(n-1)$ encrypted, a customer (a active user a) communicates with b th parties. The performance data exchange through a with privacy. The computation that describes the ratings of the target item (the item that a is looking for a prediction) finds prediction and tests a . Since data exchange are required here either a customer asks a prediction and either part can act as a active user in multiplicity scenario. Therefore, although other party's data, the proposed practical should be secure against such attacks coming from both parties.

Current proposed schemes for PPCF [1,2]. A computation of users can compute personalized recommendation with utility preserving individual data using such schemes. Patel and Du used randomized perturbation techniques for PPCF [6,7]. Vaid and Clifton [8,9,10] present privacy-preserving methods for association rule mining, χ^2 -based association classifier, and K -means clustering based on VPD. We used the CF algorithm proposed by [4]. If v_{ij} is user i 's rate for item j , and \bar{v}_i and σ_i are the mean and the standard deviation of the user i 's ratings, respectively, then the z-scores (z_{ij}) can be defined as $z_{ij} = (v_{ij} - \bar{v}_i) / \sigma_i$. Here, check the final prediction as follows where n is the number of users:

$$p_{aq} = \bar{v}_a + \sigma_a \cdot \frac{\sum_{i=1}^n w_{ai} \cdot z_{iq}}{\sum_{i=1}^n w_{ai}} \quad w_{ai} = \sum_k z_{ak} \cdot z_{ik} \tag{1}$$

here k is the item set both a and the user i have rated and q is the target item. σ_a and σ_i are standard deviations of a 's ratings and i 's ratings, respectively. p_{aq} is the prediction for a on q and w_{ai} is similarity between a and i . We used homomorphic property for our proposed practical: $E_k(x) * E_k(y) = E_k(x + y)$. Many such systems exist, and the simplest is the system by Paillier [5]. A useful property of homomorphic encryption schemes is that an addition operation can be conducted based on the encrypted data without decrypting them.

2 PPCF on VPD

With utility as a concern, data users exchange their data to provide CF services. However, with privacy as a concern, the computation should not be able

to hear each other's data. We can rewrite Eq. (1) as $p_{aq} = \overline{v}_a + \sigma_a \cdot P$ where P can be defined as follows:

$$P = \frac{\sum_{i=1}^n \left[\sum_k z_{ak} z_{ik} \right] z_{iq}}{\sum_{i=1}^n \sum_k z_{ak} z_{ik}} = \frac{\sum_k z_{ak} \left[\sum_{i=1}^n z_{ik} z_{iq} \right]}{\sum_k z_{ak} \left[\sum_{i=1}^n z_{ik} \right]} \tag{2}$$

CF systems can be either a or q oriented, rather than being much her/she oriented. Thus, p_{aq} is compared with a threshold (τ). If $p_{aq} \geq \tau$, q is recommended as like, otherwise it is recommended as dislike. If the ratings are from 1 to 5, τ is set to 3.5 (here it is set to 2 if the range from -10 to 10). Since A 's and B 's data is used to calculate P , Eq. (2) can be written as:

$$P = \frac{\sum_{k_A} z_{ak_A} \left[\sum_{i=1}^n z_{ik_A} z_{iq} \right] + \sum_{k_B} z_{ak_B} \left[\sum_{i=1}^n z_{ik_B} z_{iq} \right]}{\sum_{k_A} z_{ak_A} \left[\sum_{i=1}^n z_{ik_A} \right] + \sum_{k_B} z_{ak_B} \left[\sum_{i=1}^n z_{ik_B} \right]} = \frac{A_N + B_N}{A_D + B_D} \tag{3}$$

where $k = k_A + k_B$, and k_A and k_B represent the item sets both a and i have rated among the items held by A and B , respectively.

2.1 Off-Line Computation

The denominator part of Eq. (3) can be easily computed because A and B can find A_D and B_D using their own data. However, the numerator depends on q needs to have z_{iq} values for $i = 1, \dots, n$ to compute $\sum_{i=1}^n z_{ik_j} z_{iq}$. Efficiently for the numerator. Since A and B follow the same steps, we will explain the procedure for A . Although it divides its $n \times m_A$ data matrix into c_A submatrices where each sub-matrix consists of n/c_A users and their ratings for items held by A , where m_A is the number of items A has. It then disguises data in each sub-matrix independently. For $i = 1, \dots, c_A$, A performs the following steps:

- Step 1.** Permutes m_A columns using a permutation function Π_{Ai} .
- Step 2.** For $j = 1, \dots, m_A$, divides the permuted column $\Pi_{Ai}(I_{ij})$ into d_{ij} random vectors where $\Pi_{Ai}(I_{ij}) = \sum_{z=1}^{d_{ij}} X_{ijz}$ and d_{ij} is a 1-teger chosen with uniform random distribution over the range $[1, \beta_A]$.
- Step 3.** Permutes $X_{i11}, X_{i12}, \dots, X_{i1d_{i1}}, X_{i21}, X_{i22}, \dots, X_{i2d_{i2}}, \dots, X_{im_A1}, X_{im_A2}, \dots, X_{im_Ad_{im_A}}$ random vectors found in step 2 using π_{Ai} .
- Step 4.** Adds D_{Ai} permuted random vectors to B where $D_{Ai} = d_{i1} + d_{i2} + \dots + d_{im_A}$. B computes the scalar products between these permuted random vectors and its m_B columns using the corresponding parts of them and finds $D_{Ai}m_B$ scalar product results.
- Step 5.** Becreates the scalar product results using a homomorphic encryption scheme and its public key e_b and sends $D_{Ai}m_B$ encrypted values to A .

Step 6. Since A knows Π_{A_i} and π_{A_i} and homomorphic encryption is used, it finds the scalar product results of its m_A and B 's m_B cumulative vectors. Encrypted forms using homomorphic encryption are perturbed. After conducting these steps for $i = 1, \dots, c_A$, A gets encrypted scalar product results for its c_A sub-matrices. Since A 's data is homomorphically encrypted, it again uses homomorphic encryption to find the final scalar product results. Encrypted forms.

A creates a matrix Σ_A consisting of $e_b(\Sigma_{ij})$ for $i = 1, \dots, m_A$ and $j = 1, \dots, m_B$ here $e_b(\Sigma_{ij})$ represents the encrypted scalar product between i^{th} cumulative vector of B and the j^{th} cumulative vector of A . It generates large enough v_{ij} random numbers for $i = 1, \dots, m_B$ and $j = 1, \dots, m_A$, encrypts them using e_b , and adds them to the $e_b(\Sigma_{ij})$ values using homomorphic encryption. It finds matrix Σ'_A consisting of $e_b(\Sigma'_{ij})$ here $\Sigma'_{ij} = \Sigma_{ij} + v_{ij}$ and stores v_{ij} values in a matrix V_A . It sends Σ'_A to B that decrypts the encrypted values and finds the matrix Σ''_A consisting of Σ'_{ij} values and stores it. Before doing the same procedure, B finds matrices Σ_B, Σ'_B , and V_B . B stores V_B and finds Σ''_B , and then stores it. A and B compute the item means and store them in $m_A \times 1$ and $m_B \times 1$ matrices, respectively.

2.2 Online Computation

Since either party can act as a active user in multiparty scenarios, the computation should be secure against such attacks. The steps are as follows:

Step 1. a sends his/her data and a query to the compiler that sends q . Assume that B sends q . B computes B_N and B_D . However, since A can act as a active user in multiparty scenarios, therefore, B uses private B_N & B_D computation protocol, which is explained in the following, to compute them.

Step 2. B can compute A'_N value using the data from the q^{th} row of the matrix Σ''_A and a 's corresponding data here $A'_N = A_N + R_q$. The data from the q^{th} row of the matrix Σ''_A represents $\sum_{i=1}^n z_{ik_A} z_{iq}$ values disguised by v_{qk_A} random numbers for k_A . A can compute $R_q = \sum_{k_A} z_{ak_A} v_{qk_A}$ here k_A represents the items rated by a among the items held by A .

Step 3. B computes $A_N + R_q + B'_N$ and B'_D and sends them together with a 's item mean, standard deviation, and the z-scores for the items rated by a among q items held by A to A through a . A computes R_q , finds $A_N + B'_N = A'_N + B'_N - R_q$ and A_D , and estimates P' using Eq. 3 based on the query.

A computes p'_{aq} , then a whether he/she is like q or not by comparing p'_{aq} with τ . Since B can act as a active user, A uses a random threshold to prevent B from learning A_D and A_N . It generates a uniform random number $(r_{A\tau})$ from a range $[-\alpha_A, \alpha_A]$, finds $\tau + r_{A\tau}$, and uses it as a random threshold.

Our scheme can be extended to multiparty. Each edge receives data from $n-1$ entities and stores it as a n -party scheme. During the phase, a sends his/her data to the party that sends q . That party computes the required data like it does in n -party scheme and sends results to a . The n th party acts as a master site. Other parties compute the values required for enumeration and deliver them to all parts. Each compiler creates a large enough uniform

ber from a range $[-\gamma, \gamma]$, adds it to the values for numerator and denominator parts, and sends them through a trusted master site, which estimates the prediction.

Private B_N & B_D Computation Protocol. We explain the protocol for B . After B gets a' data, it finds the number of rated items (C_B) that a rated among the items it holds. If C_B is less than $\lfloor m_B/2 \rfloor$, then B finds the items that a did not rate among the items B holds. B generates a uniform random integer S_{Ba} from the range $[1, m_B - C_B]$, random selects S_{Ba} unrated items among the items it holds, and finds their average as the a 's ratings reflect with their mean votes. If C_B is bigger than $\lfloor m_B/2 \rfloor$, B finds the items that a rated among the items B holds and creates a uniform random integer S_{Br} from the range $[1, C_B]$. It random selects S_{Br} rated items and removes their ratings from a 's ratings reflect. B computes B'_N and B'_D using the new ratings reflect of a and finds a ratings' average and standard deviation and computes the z-score.

3 Privacy and Overhead Costs Analysis

In this section we first investigate privacy.

Claim 1. B learns nothing about A 's data. A 's data is $(m_A!D_A!(\beta_A)^{m_A})^{c_A}$. Since A uses Π_{A_i} s for $i = 1, \dots, c_A$ to permute its m_A columns reflectors in each sub-matrix for B , the probability of guessing the correct position of them is 1 out of $m_A!$. A decides each of its permuted columns in the random reflectors here it decides how many random reflectors a permuted reflector be divided into based on a uniform random integer from the range $[1, \beta_A]$. Therefore, the probability of guessing the number of random reflectors that each reflector is divided into is 1 out of $(\beta_A)^{m_A}$. A uses π_{A_i} s for $i = 1, \dots, c_A$ to permute random reflectors. Therefore, guessing their correct positions is 1 out of $D_A!$ with the assumption that all D_{A_i} values are same and equal to D_A . Since A hides its data into c_A parts, the probability of guessing the A 's data for B is 1 out of $(m_A!D_A!(\beta_A)^{m_A})^{c_A}$.

Claim 2. A learns nothing about B_N & B_D . B_N & B_D are $(m_B - C_B) \times C_B$. Since B uses random integer S_{Ba} and random selects S_{Ba} unrated items among the items it holds, the probability of guessing the correct S_{Ba} and which S_{Ba} unrated items are selected is 1 out of $((m_B - C_B)(m_B - C_B)!)/(S_{Ba}!(m_B - C_B - S_{Ba}!))$. B also finds unrated items' average as a 's ratings reflect with their mean votes, which are unknown to B .

Claim 3. B learns nothing about A_N & A_D . A_N & A_D are $(m_A - C_A) \times C_A$. Since A tests that he/she is like or dislike q and produces referrals using a random threshold, B cannot hear A_N and A_D .

Claim 4. A learns nothing about B . B is $(m_B - C_B) \times C_B$. B 's data is $\Sigma_{ij} \dots \Sigma''_{ij} \dots$.

Unlike the previous communication cost, the communication cost is instead the number of the communication is 4 for our scheme. The additional storage costs due to privacy issues are $O(m_A m_B + m_A)$ and $O(m_A m_B + m_B)$ for

A and B , respectively. Although the time complexity is not critical, the space complexity is essential.

Claim 5. $O(c_B D_B + m_A m_B)$ for A and $O(c_A D_A + m_A m_B)$ for B .

Claim 6. $O(c_B D_B + m_A m_B)$ for A and $O(c_A D_A + m_A m_B)$ for B , where D_A and D_B are the number of items rated by A and B , respectively.

Claim 7. $O(nm_B D_A)$ for A and $O(nm_A D_B)$ for B .

4 Experiments

We used Jester and MovieLens (ML) data sets. Jester [3] has 100 jokes and records of 17,988 users. The ratings range from -10 to +10. ML (www.cs.um.edu/research/Group) consists of ratings for 3,592 movies made by 7,463 users. Ratings are made on a 5-star scale. We measured the accuracy of our approach using classification accuracy (CA) and F -measure (FM), which is a weighted combination of precision and recall.

4.1 Methodology

We randomly selected 2,000 users for training from Jester and ML. Since we conducted different sets of experiments with varying number of rated items (M), we used these users who rated certain number of items and randomly selected 400 test users among them for each experiment. For each test user, we randomly selected 5 rated items, including a single rated item for each test user, and tried to predict its true genre and other ratings. We did this for a 5 test items. We replaced the test item's entries as usual. We randomly selected a subset of rated or unrated items proportionally 10 times for each test item. We created r_{τ} uniform random numbers higher than the rating evaluation threshold. For each test item, we created 10 uniform random numbers for these experiments testing accuracy with random threshold. We converted the threshold items' ratings into binary ratings. We then compared the recommended results of our scheme with the threshold items' converted ratings and found CA and FM values.

4.2 Experimental Results

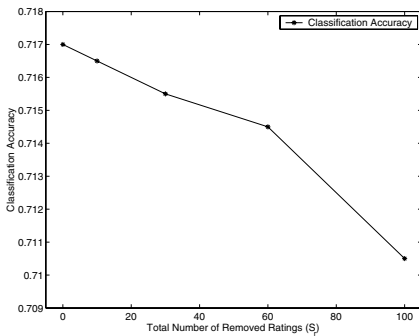
Number of Rated Items (M). We hypothesize that since prediction quality improves with increasing M , the test parties conduct CF on the initial data, accuracy improves. To show the effects of different M values, we conducted experiments using ML data with varying M . Table 1 shows CAs and FMs with varying M . Based on the settings of each experiment, we selected these users for testing who rated M number of items. Overall performance increases with increasing M . If there is limited number of rated items, with increasing M values, we gain significant improvement. However, the improvement becomes

Table 1. Prediction Quality vs. M

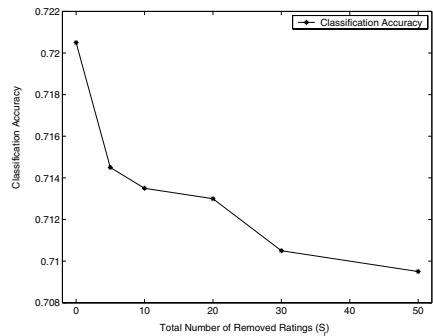
M	$M < 50$	$40 < M < 100$	$100 \leq M < 200$	$200 \leq M < 400$
CA	0.6645	0.7010	0.7100	0.7140
FM	0.7313	0.7686	0.7713	0.7743

stable enough ratings are available. CA is 0.6645 when M is less than 50. Hence it increases to 0.7010 when $40 < M < 100$. Besides CA, FM also increases from 0.7313 when $M < 50$ to 0.7686 when $40 < M < 100$.

Number of Removed Ratings (S_r). We conducted experiments where varying S_r using Jester and ML, and showed CAs in Fig. 1. 400 test users were randomly selected among these users who rated more than 200 and 80 items for ML and Jester, respectively. As seen from Fig. 1, accuracy slightly becomes worse with increasing S_r because the available ratings are decreasing. When we increased S_r from 0 to 100, the best 0.0065 accuracy for ML. This means that if there are significant large number of ratings available, removing some of them does not affect accuracy too much.



(a) ML (rating range: 1-5)



(b) Jester (rating range: -10-10)

Fig. 1. Prediction Quality vs. S_r

Number of Appended Ratings (S_a). Since accuracy improves with increasing available ratings, appending more ratings may improve accuracy. However, since empty cells are filled with the item mean values, which can be considered default values for a random match with his/her true ratings for these items, that might make accuracy worse. We performed experiments using ML with varying S_a . 400 test users were randomly selected among these users who rated more than 40 and less than 100 items. CA improves with increasing S_a up to 60 appended ratings. When S_a is 100, CA becomes worse and it is 0.7035; but it is still better than the CA, which is 0.7010, when S_a is 0.

Range of Uniform Random Values (α). The shift in direct α values affects our results, error percentages highly depend on the created r_τ values from the range $[-\alpha, \alpha]$ with varying α values using ML. 400 test users were randomly selected among these users. We noted more than 40 and less than 100 items. The results slightly become worse with increasing α because with increasing range, r_τ values become larger, the random threshold $(\tau + r_\tau)$ fluctuates more and causes a loss in the performance. When we increased α from 0 to 0.1, CA degrades by 0.0015 while FM decreases by 0.0032.

5 Conclusions and Future Work

We have presented a solution to the PPCF based VPD problem. Our solution makes it possible for the parties to conduct their services using their own data without disclosing their data to each other. Our experiment results have shown that our solution produces accurate referrals compared with the true ratings. We also studied multi-part scheme in detail and showed its accuracy and practical change compared to the k -part scheme and with varying number of parties.

References

1. J. Canny. Collaborative filtering with privacy. In *Proceedings of the IEEE Symposium on Security and Privacy*, pages 45–57, Oakland, CA, USA, May 2002.
2. J. Canny. Collaborative filtering with privacy via factor analysis. In *Proceedings of the 25th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 238–245, Tampere, Finland, 2002.
3. D. Gupta, M. Digiovanni, H. Narita, and K. Goldberg. Jester 2.0: A new linear-time collaborative filtering algorithm applied to jokes. In *Workshop on Recommender Systems Algorithms and Evaluation, 22nd International Conference on Research and Development in Information Retrieval*, Berkeley, CA, USA, August 1999.
4. J. L. Herlocker, J. A. Konstan, A. Borchers, and J. T. Riedl. An algorithmic framework for performing collaborative filtering. In *Proceedings of the 1999 Conference on Research and Development in Information Retrieval*, August 1999.
5. P. Paillier. Public-key cryptosystems based on composite degree residue classes. In *Advances in Cryptology – EUROCRYPT’99*, pages 223–238, 1999.
6. H. Polat and W. Du. Privacy-preserving collaborative filtering using randomized perturbation techniques. In *Proceedings of the 3rd IEEE International Conference on Data Mining (ICDM’03)*, Melbourne, FL, USA, November 19–22 2003.
7. H. Polat and W. Du. SVD-based collaborative filtering with privacy. In *Proceedings of the 20th ACM Symposium on Applied Computing Special Track on E-commerce Technologies*, Santa Fe, NM, USA, March 13–17 2005.
8. J. Vaidya and C. Clifton. Privacy preserving association rule mining in vertically partitioned data. In *Proceedings of the 8th ACM SIGKDD*, pages 639–644, 2002.
9. J. Vaidya and C. Clifton. Privacy preserving k-means clustering over vertically partitioned data. In *Proceedings of the 2003 ACM SIGKDD*, Washington, DC, USA, August 24–27 2003.
10. J. Vaidya and C. Clifton. Privacy preserving naïve bayes classifier for vertically partitioned data. In *Proceedings of the 2004 SIAM Conference on Data Mining*, Orlando, FL, USA, May 2004.

Indexed Bit Map (IBM) for Mining Frequent Sequences

Lionel Savary and Karine Zeitouni

PRiSM Laboratory, 45 Avenue des Etats-Unis,
78035 Versailles, France
{Lionel.Savary, Karine.Zeitouni}@prism.uvsq.fr

Abstract. Sequential pattern mining has been an emerging problem in data mining. In this paper, we propose a new algorithm for mining frequent sequences. It processes only one scan of the database thanks to an indexed structure associated to a bit map representation. Thus, it allows a fast data access and a compact storage in main memory. The experimental results show the efficiency of our method compared to existing algorithms. It has been tested on synthetic data and on real data containing sequences of activities of a urban population time-use survey.

1 Introduction

The problem of mining sequential patterns was first introduced in the context of customer transactions analysis [2]. It aims to retrieve frequent patterns in the sequences of products purchased by customers through time ordered transactions. Several algorithms have been proposed in order to improve the performances and to reduce required space in memory [5], [9], [6]. Other works have concerned mining frequent sequences in DNA [8] or Web Usage Mining [3]. Finally, notice the use of bit map structure in providing a compact representation and good performances [5].

The target application in this paper is related to population time-use analysis and more precisely their daily displacements [4]. Our data are related to daily activities carried out by each surveyed person at the scale of a whole urban area. Thus, for each person of a surveyed household, it captures the activity program [7], the transport mode used between two activities, the departure time, and the duration of the trip. For example, during a day, an individual can leave home, take children to school, go to work, pick children up from school, and come back home. Activity programs of most individuals may be the same or be similar. Each activity program could be seen as a sequence of single values, making it possible to discover frequent activity sequences that characterise groups of the surveyed individuals. This allows analyzing the mobility of this urban population. Likewise, when considering transport mode, schedules or duration sequences, it would be possible to determine a typology of used transport modes, schedules, and so on.

Existing algorithms are either inappropriate or not enough efficient to our specific case. Most works [1], [2], [6] make multiple scan of the database, which can be considered as the main bottleneck of algorithms of frequent sequence mining. Furthermore, unlike the analysis of sequential transactions where each transaction is an item set, our context only focuses on the analysis of sequences of items.

Although existing works [9], [10], [12] can be applied in this context, we propose here a new algorithm more appropriate to this particular case. This algorithm only makes one scan of the database. The indexed bit map structure needs few spaces in the main memory and allows a fast access to the data. The experimental results, using real or synthetic data, show that our algorithm outperforms existing ones.

The paper is organised as follows: section 2 presents related works, then, section 3 describes the proposed algorithm, section 4 proposes an optimisation, section 5 relates the experimentation and performance study, and finally, a general conclusion summarizes our contribution and traces some perspectives.

2 Related Works

Most works related to mining frequent sequences are in the field of customer transaction analysis. Early work on frequent patterns -*Apriori* algorithm- only considered transactions, not sequence of transactions [1]. This algorithm is costly because it carries out multiple scans of the database to determine frequent subsets of items. Three algorithms dealing with sequence of transactions are presented and compared in [2]: *AprioriAll*, *AprioriSome* and *DynamicSome*. *AprioriAll* algorithm is an adaptation of *Apriori* to sequences where candidate generation and support are computed differently. *AprioriAll*, and *AprioriSome* only compute maximal frequent sequences. Their principle is to jump to candidates of size $k+next(k)$ in the next scan, where $next(k)>1$. Maximum frequent sequences of lower size that have not been calculated are given in the backward phase. The value of $next(k)$ increases with $P_k = |L_k|/|C_k|$, where L_k stands for frequent sequences of size k , and C_k the whole generated candidates of size k . *DynamicSome* algorithm is based on *AprioriSome* but uses a jump by a multiple of user defined *step*.

SPAM algorithm [5] uses a bitmap representation of transaction sequences once the entire database has been loaded in a lexicographic tree. But this algorithm considers that the entire database and all used data structures should completely fit into main memory, and then do not adapt for large datasets.

The *GSP* algorithm [6] exploits the property that all contiguous subsequences of a frequent sequence also have to be frequent. As *Apriori*, it generates frequent sequences, then candidate sequences by adding one or more items.

PrefixSpan [10] first finds the frequent items after scanning the database once. The sequence database is then projected, according to the frequent items, into several smaller databases. Finally, all sequential patterns are found by recursively growing subsequence fragments in each projected database. Employing a divide-and-conquer strategy with the *PatternGrowth* methodology, *PrefixSpan* efficiently mines the complete set of patterns.

3 IBM Algorithm

We are now going to focus on the specific case where the considered sequences are basic since they are composed of single items, not of a set of items. This is the case in DNA [8], Web usage data [3] or activity program sequences [7]. Our algorithm will

be compared to PrefixSpan, one of the most efficient among the above mentioned methods.

A sequence is said frequent if it is included in a number of sequences greater than a support given by the user. The inclusion between two sequences $s_1 = (a_1, \dots, a_n)$ and $s_2 = (b_1, \dots, b_n)$: $s_1 \subset s_2$ is defined by : $\exists b_{i_1} = a_1, \dots, b_{i_n} = a_n$ such that $i_1 < i_2 < \dots < i_n$.

3.1 Principle of the Algorithm

The proposed approach is two phases. The first stage is the data encoding into a memory resident data structures. The second one is the frequent generation that in turn is composed of candidate generation, and candidate support checking.

The data structure is based on four components: (i) a Bit Map (IBM) is a binary matrix representing the distinct sequences of the database, (ii) an SV vector encodes all the ordered combinations of sequences, (iii) an index (INDEX) on the Bit Map allows a direct access to sequences according to their size, (iv) an NB table associated to the Bit Map which informs about the frequency of each distinct sequence (figure 1).

This algorithm only makes one scan of the database during which the total number of distinct sequences, the frequency of these sequences and the number of sequence by size are computed. This allows computing the support of each generated sequence. These sequences are classified by decreasing size in the IBM and only distinct sequences are stored in the Bit Map. An index by size allows a direct access to sequences according to their size. This structure provides an optimisation since a generated sequence s of size t will be directly compared with the sequences of the same or greater size stored in the IBM (figure 1).

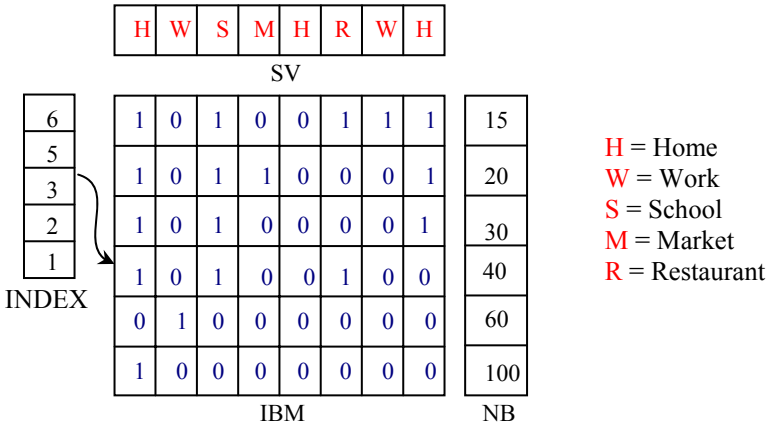


Fig. 1. The data structure

In order to simplify the notations, we represent each activity by a specific character, e.g. HSWSH standing for (Home, School, Work, School, Home). In the figure 1, the sequence vector (SV) is made of 5 ordered activities (H,W,S,M,H). In this example one supposes that the database is composed of six distinct sequences of size 1 to 5

encoded in the IBM. The bit 1 indicates the items present in the sequence according to the SV and bit 0, those that are not. Here, there are 6 distinct sequences: (H), (W), (HRS), (HSH), (HSMH) and (HSRWH). In the above example (figure 1), each cell of the INDEX indicates the first line where the corresponding size of sequence is stored. For example, the cell number 5 (with value 6) corresponds to the line number 6 of the first sequence of size 5 encoded in the IBM. The table NB associates to the IBM stores the frequency of each distinct sequence. Thus the sequence (HSMH) of size 4 occurs 20 times in the database. In this algorithm, INDEX, SV, NB and IBM are built on the fly during one pass. At each insertion of a sequence, the IBM may become larger, and a set of shifting operations are applied to the bit values stored in this table.

```

IBM (sequence database DB, threshold t)
00 For each sequence s in DB
01   Gen-sequence-vector(s)
02   Encode and Insert s in the IBM
03   Update NB
04   Update INDEX
05 End For
06 Integer k := 1;
07 While exists frequent sequence of size k
08   k := k+1;
09   Generate Ck
10   Get-frequent-sequences (t)
11 End While

```

Fig. 2. IBM algorithm

Figure 2 shows the general IBM algorithm that takes as parameters: the database of sequences DB and a threshold t . This value (t) stands for the minimum frequency of the sequences which will be taken into account for the generation of the candidates. Then for each sequence s reads from the database during the scan, the SV (line 01) is generated using a merging process (see section 3.2). If the sequence already exists in SV, only the NB table is updated (line 03): the line corresponding to this sequence in NB (and encoded in the IBM) is incremented. So, the frequency corresponding to this value is incremented. Else, if the sequence is not presented in SV, it is generated by the Gen-sequence-vector(s) function (section 3.2). The height of the IBM is increased to one line (line 02), the length is increased to the SV length, and the INDEX (line 04) is updated. Then, a set of shifting operations is applied to the IBM in order to preserve the initial values of existing sequences while encoding the new one.

Once all the data have been encoded in this structure (SV, IBM, NB, INDEX), new candidates (line 09) are generated (see section 3.3) and compared to the data stored in the IBM (line 10) with a fast access thanks to the index (INDEX).

3.2 Generation of the Sequence Vector

The sequence vector is generated during the unique scan of the database according to the algorithm of figure 3. Here, s stands for a sequence of the database read during the

scan, and $\text{position}(x)$ stands for the cell number of value x in the SV. If an item a of s already exists in SV, then there is nothing to do, otherwise, there are two possibilities: if there exists an item b such that the cell number of b is greater than the cell number of a and b is in SV (line 04 and 05), then a is inserted before the value b in SV; otherwise, a is inserted at the end of SV (line 06). Thus all the distinct sequences of the database are represented in the SV using a merging process.

```

Gen-sequence-vector( $s$ ):
00 var SV :=  $\emptyset$ ; {SV empty at the beginning};
01 Integer current_position := 0; {position in SV};
02 For each item  $a$  of  $s$ 
03   If  $a \notin$  SV
04     If  $\exists b \in s$  such that ( $b \in$  SV and  $\text{position}(b) >$ 
 $\text{position}(a)$  in  $s$  and  $\text{position}(b) >$  current_position)
05       Insert  $a$  before  $b$ 
06     Else insert  $a$  at the end of SV
07       current_position :=  $\text{position}(a)$  in SV;
08 End For

```

Fig. 3. Sequence Vector generation

3.3 Candidate Generation

During the scan, the frequencies of all items are computed. Those whose support is underneath the one specified by the user are deleted. Then, candidates are generated from these frequent items, using the fusion process as in GSP algorithm [6].

3.4 Candidate Support Counting

For a given candidate C of size S , the algorithm first accesses the first sequence of size S encoded in IBM, which corresponds to the line $l = \text{INDEX}(S)$. For each line starting from the line l to the last line of IBM table, the algorithm determines using the SV vector if C is contained in each line of IBM. If so, the corresponding frequency of this sequence stored in the NB table, is added to the frequency of the candidate. After the comparison with each line until the last one, the support of C is computed.

4 Implementation and Optimization

The IBM algorithm has been implemented in Java. It takes few spaces in the main memory. But whereas the bit variable is not provided in programming languages like Java or C++, some shifting operations are required to access the target value stored in the bit map and corresponding to the value stored in SV. In order to avoid these superfluous computations, we have proposed a variant with IBM2 algorithm, where the bit map is replaced by a Boolean matrix, i.e. where cells are declared of Boolean type, which takes 8 bits for each cell. Although this solution requires more space in mem-

ory, the access to the target value stored in the Boolean matrix is done directly without shifting computations. The result of their respective performances is detailed in the next section and compared with PrefixSpan.

5 Experimental Results

The experiments were performed on a 2.5 GHz Pentium IV with 1.5 GB of memory running Microsoft Windows XP Professional. Our implementation of IBM and IBM2 has been compared with PrefixSpan, based on the package PrefixSpan-0.4.tar.gz¹. This test has concerned the scalability of the algorithm, by measurements of runtime and memory occupancy while varying the dataset size, and the support threshold. Moreover, we have tested the impact of the number of distinct items. Four synthetic datasets have been generated for the experimentations, with different sizes: 100000, 300000, 600000 and 1000000 rows. The size of sequences is randomly generated from 2 to 60, and the number of distinct items is about 10 for figures 4 to 7. This number has been pushed to 35 distinct items in order to test its impact.

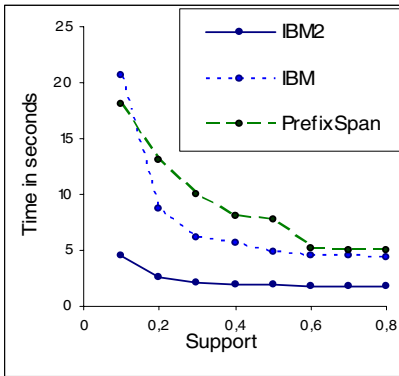


Fig. 4. Performances with 100,000 rows

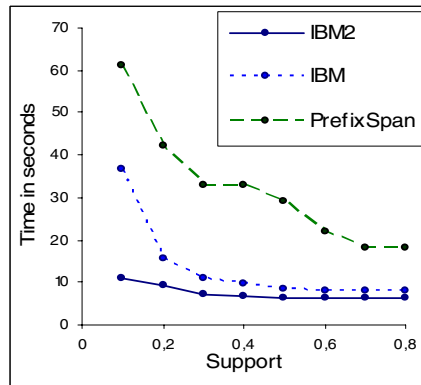


Fig. 5. Performances with 300,000 rows

Although IBM and IBM2 have been implemented in Java, and PrefixSpan in C++ - *a priori* more optimal than Java -, IBM and IBM2 outperform PrefixSpan. The experimentations show that the larger is the database size, the more IBM and IBM2 win PrefixSpan (Figures 4 to 7). This is because IBM and IBM2 make only one scan of the database and the Indexed Bit Map structure allows a faster access to the sequences than the structure used in PrefixSpan. Moreover, as the support threshold decreases, the gap between IBM and PrefixSpan increases. Concerning the resource consumption, the size of the bit map depends on the size of SV, which may increase with the number of distinct sequences. Notice that SV size only increases when the encountered sequence can not be encoded using the current SV. Moreover, not all the items of the inserted sequence are added in SV, but only those that are not present in the

¹ <http://chasen.org/~taku/software/prefixspan/>

same order. Finally, since the probability to find common ordered items between SV and the current sequence becomes high as the building process advances, SV size becomes stable regardless of the size of the database. For instance, with a database composed of 600,000 rows, SV contains about 265 values for 90,000 distinct rows. The size of the Boolean Map is then equal to: $265 \times 90,000 = 23.85$ Mega Bytes. As IBM is 8 times more compact, the size of the Bit Map is less than 3 MB. With 1,000,000 rows (figure 7), SV contains 370 elements for 160,000 distinct rows. Then, the size of the Boolean Map reaches 59.2 MB, whereas the size of the Bit Map fits in 7.5 MB. Concerning the impact of distinct item number, for 100,000 rows until 20 distinct items, IBM and IBM2 perform better than PrefixSpan. Between 20 and 35 distinct items, IBM2 performs better than PrefixSpan, which becomes faster than IBM. But above 35 distinct items, PrefixSpan is faster than IBM and IBM2.

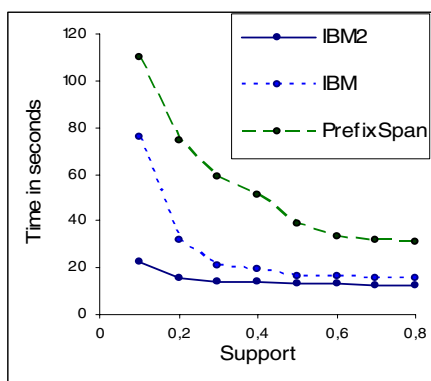


Fig. 6. Performances with 600,000 rows

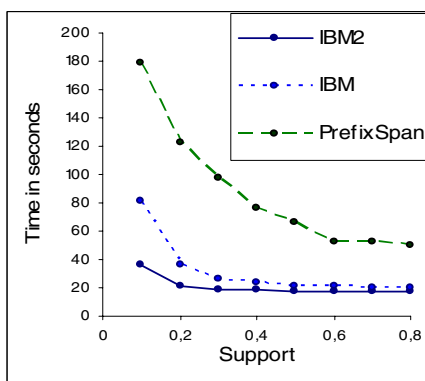


Fig. 7. Performances with 1,000,000 rows

These results also show that IBM is more appropriate than IBM2 for very large databases, due to data compression. However, IBM2 runs faster than IBM. This is due to the costs of shifting operations necessary to access target values, while IBM2 directly accesses the target sequences.

6 Conclusion and Perspectives

This paper has presented a new algorithm IBM and its variant IBM2. The aim of this algorithm is to find all frequent sequences in item sequences. It has been applied to discover all frequent activity sequences in the time use mobility database within an urban environment. IBM only makes one scan of the database and provides an efficient data structure saving runtime and memory space. The use of the specified index provides another optimization of comparisons during candidate counting. Experimental results show that in most cases, IBM2 outperforms IBM, which in turn outperforms PrefixSpan for large and very large databases, with limited distinct items. Extensive experiments have been conducted that attest for the effectiveness and the efficiency of the proposed method, and are detailed in [11]. In perspective, IBM will

be extended to multidimensional sequences (e.g. with attributes) and spatial sequences (such as trajectories). Other application fields will be explored, like pattern mining from DNA, Web Usage Mining or extension to customer transaction analysis. Finally, the proposed data structure adapts to similarity analysis of sequences and may be a good basis for efficient sequence clustering.

References

1. Agrawal, R., Srikant, R.: Fast Algorithms for Mining Association Rules. In Proc. of the 20th Int. Conf. Very Large Data Bases (VLDB), Santiago, Chile, September (1994)
2. Agrawal, R., Srikant, R.: Mining sequential patterns. In Proc. of the 11th Int'l Conference on Data Engineering, Taipei, Taiwan, March (1995)
3. Spiliopoulou, M., Faulstich Lukas, C., Winkler, K.: A data miner analyzing the navigational behaviour of web users. In Proc. Of the Workshop on Machine Learning in User Modelling of the ACAI'99 Int. Conf., Crete, Greece, July (1999)
4. Ministère de l'Équipement, des Transports et du Logement. L'enquête ménages déplacements « méthode standard ». Collections du Certu. Octobre (1998), ISSN 1263-3313
5. Jay, A., Johannes, G., Tomi, Y., Jason F.: Sequential Pattern Mining using A Bitmap Representation. SIGMOD pp 429-435, July (2002), Edmonton, Alberta, Canada
6. Srikant, R., Agrawal, R.: Mining Sequential Patterns : Generalizations and Performance Improvements. Proc. 5th EDBT, Mars 25-29, (1996). Avignon, France. pp 3-17
7. Wang, D., Tao, C.: A spatio-temporal data model for activity-based transport demand modeling.. International Journal of Geographical Information Science, (2001), 15(6), pp 561-585
8. Han, J., Jamil, H. M., Lu, Y., Chen, L., Liao, Y., Pei, J.: DNA Miner: A system prototype for mining DNA sequences. In the proc. Of the ACM SIGMOD International Conference on the management of data, Day 21-24,(2001), Santa Barbara, CA, USA
9. Zaki, M. J.: Efficient Enumeration of Frequent Sequences. Int. Conference on Information and Knowledge Management, November(1998), Washington DC
10. Pei, J., Han, J., Mortazavi-Asl, B., and H., Pinto.: Prefixspan: Mining sequential patterns efficiency by prefix-projected pattern growth. In Proc. of the International Conference on Data Engineering (ICDE), pp 215–224, (2001)
11. Savary, L., Zeitouni, K.: Indexed Bit Map (IBM) for Mining Frequent Sequences. PRiSM Laboratory Technical Report N° 2005/82, August (2005), Versailles University, France. <http://www.prism.uvsq.fr/rapports/bin/bibliography.php?id=300>

STochFS: A Framework for Combining Feature Selection Outcomes Through a Stochastic Process

Jefferson Teixeira de Sá¹, Nathaniel Jacobson², and Samuel M. Elomaa^{1,2}

¹ Computer Science Department, Federal University of Ceará,
Fortaleza, 60455-760, Brazil
jeff@lia.ufc.br

² School of Information Technology and Engineering, University of Ottawa,
Ottawa, K1N 6N5, Canada
{nat,stan}@site.uottawa.ca

Abstract. The *Feature Selection* problem involves discovering a subset of features such that a classifier built only with this subset would have better predictive accuracy than a classifier built from the entire set of features. Ensemble methods, such as Bagging and Boosting, have been shown to increase the performance of classifiers to remarkable levels but surprisingly have not been tried in other parts of the classification process. In this paper, we apply the ensemble approach to feature selection by proposing a systematic way of combining various outcomes of a feature selection algorithm. The proposed framework, named STochFS, have been shown empirically to improve the performance of well-known feature selection algorithms.

1 Introduction

The feature selection problem involves discovering a subset of features such that a classifier built only with this subset would have better predictive accuracy than a classifier built from the entire set of features.

Ensemble methods, such as Bagging and Boosting, have been shown to increase the performance of classifiers to remarkable levels but surprisingly have not been tried in other parts of the classification process. In this paper, we apply the ensemble approach to feature selection by proposing a systematic way of combining various outcomes of a feature selection algorithm. The proposed framework, named STochFS, have been shown empirically to improve the performance of well-known feature selection algorithms.

2 Feature Selection

Feature selection algorithms can be categorized into two main groups: filter methods and wrapper methods.

ea, i g a g, i h ed c c he ca i e. If fea e eec i i e- f ed i de e de f he ea, i g a g, i h, he ech i e i aid f a a a ach. O he i e, i i aid f a a a ach. Whe he e a a ach i ge e a c a i a e e e cie ha he a e a a ach, i a a da bac i ha a i a eec i f fea e a be i de e de f he i d c i e a d e e e a i a bia e f he ea, i g a g, i h ha i ed c c he ca i e. The a e a ach he he ha d, i e he c a i a e head fea a i g ca dida e fea e be b e e c i g a e e c ed ea, i g a g, i h he da a e e e e d i g each fea e be de c i de a i .

A c b i a i f he e a a che, ha i, he e f e a a i e h d (a e- e e a a i f c i a d a c a i e) c e a e a i . H b i d i a e c b i e he g d cha ac e i c f b h e a d a e¹. The c b i a i f a a che e f ed b h b i d fea e e c i a g, i h, h e e, a e he i c i a e a d c a be e e a i c a a i ed a a ed.

I h i a e, e e e he e a i c c b i a i f he c e f fea e e e c i a g, i h i g he Bagg i g ech i e i a a cha i c c e .

3 The Framework

The ST chFS fa e c b i e he e f a fea e e e c i a g, i h i a a cha i c a e b a i g he e c e i a i ge c e a d i g i a a e e d i he ge e a i f e fea e e e c i b e i h a e e a a e d i h a ea, i g a g, i h .

I i i a , he *NumOuts* be b e e ed b a i ge f a fea e e e c i e *fs* (he i ge e f *NumOuts* d i e e , if ch a a g, i h e e b e b e e e e c i) a e e d i a d i e i a a a , ee Fig e 1. Thi a a i he b e c d e d i a e a a , ca ed *Adam*, ha i i e he b e f i e each fea e a e a e d i he *NumOuts* be b e . Ne , ST chFS i i e a i e (*NumIter* i e) ge e a e e b e f fea e i a a cha i c a g i d e d fa h i i g *Adam* a a e e d a d e a e he i h a ea, i g e e he da a e *D*. The ge e a i f a e b e i ch ha fea e i h high a e i *Adam* ha e a b e e cha c e f b e i g e e c ed ha h e i h a e e a each i e a i . A he e d, he b e i h b e acc ac i b e e e d. If b e i e i e f acc ac, he e i h he e c a d i a i i e e d.

Each f he c e d e e d i h i f a e a e d e c i b e e .

GenerateOutcomes(*fs, D, NumOuts*) e e e he fea e e e c i a g, i h , *fs, NumOuts* i e a d e i c e i *O*. Thi c e d e e d i e e , a d e c i b e d i he e e c i , d e e d i g he he fea e e e c i a g, i h b e e d i b a b i i c , d e e i i c .

¹ A description of recently proposed hybrid feature selections algorithms can be found in [13], [14], [2],[5] and [12].

```

STochFS( $fs, D, NumIter, NumOuts$ )

 $O = \text{GenerateOutcomes}(fs, D, NumOuts)$ 
 $Adam = \text{CalculateAdam}(O)$ 
for  $j = 1$  to  $NumIter$ 
     $S = \text{GenerateSubset}(Adam)$ 
    if  $\text{Error}(S, D) < \text{Error}(S_{best}, D)$  then
         $S_{best} = S$ 
    else
        if  $\text{Error}(S, D) = \text{Error}(S_{best}, D)$  and
             $\text{Card}(S) < \text{Card}(S_{best})$  then
                 $S_{best} = S$ 
return  $S_{best}$ 

```

Fig. 1. The STochFS Framework

CalculateAdam(O) generates a vector $Adam$:

$$Adam = \{a_i, 1 \leq i \leq n\}$$

where $a_i = \sum o_{ji}$, for $1 \leq j \leq k$ and $1 \leq i \leq n$.

Calculate the $Adam$ vector. $Adam$ is the vector of frequencies of each feature in O , which is the best feature set selected by fs a given feature.

GenerateSubset($Adam$) generates a best feature S in a characteristic guided fashion using $Adam$ as a seed. The generation process is described below. Let i denote a random feature in $Adam$. Let S be a vector of frequencies of features in the current best feature set O . For S_i ($f \in S$) = 1 if feature i is included in the best feature set selected by S . $S_i = 0$, otherwise. Vector S is constructed as follows:

$S_i = 1$, if $a_i > \text{random}(k)$ and $S_i = 0$ otherwise,

where $\text{random}(k)$ is a random number between 0 and k .

This procedure is repeated for each feature with high frequencies and a best characteristic feature set is chosen. The procedure is repeated for each feature.

Error(S, D) is a performance evaluation metric, using the best feature set S generated by the procedure described above, to evaluate the performance of S on the data set D .

In the following section, we describe the generated best feature set using the STochFS framework. Section 3.1 describes the basic framework, and Section 3.2 describes the advanced framework.

3.1 Combining Outcomes of Probabilistic Feature Selection Algorithms

In the case of the feature selection algorithm used in ST chFS, fs_{11} and fs_{12} are bagging algorithms, GenerateOutcomes($fs, D, NumOuts$) proceeds as follows:

```

GenerateOutcomes( $fs, D, NumOuts$ )

  for  $i = 1$  to  $NumOuts$ 
     $O[i] = FeatureSelection(fs, D)$ 
    
```

Fig. 2. GenerateOutcomes() for Probabilistic Algorithms

here

FeatureSelection(fs, D) is a feature selection algorithm fs on D and returns the selected $O[i]$.

Since fs_{11} and fs_{12} are bagging algorithms, a different set of features is chosen in each iteration of the feature selection.

3.2 Combining Outcomes of Deterministic Feature Selection Algorithms

Otherwise, had fs_{11} been deterministic, here the GenerateOutcomes algorithm:

```

GenerateOutcomes( $fs, D, NumOuts$ )

  for  $i = 1$  to  $NumOuts$ 
     $D[i] = Resample(D)$ 
     $O[i] = FeatureSelection(D[i])$ 
    
```

Fig. 3. GenerateOutcomes() for Deterministic Algorithms

here

Resample(D) creates a new version of the original data D by sampling with replacement [4]. Each bootstrap sample $D[i]$, contains approximately the same age 63.2% of the original D^2 .

FeatureSelection($D[i]$) is a feature selection algorithm fs on $D[i]$ and returns the selected $O[i]$.

The selected features added to the deterministic feature selection are added to the set of features selected by the stochastic bagging based on *Adam* using the ST chFS.

² This is the same sampling technique used in Bagging.

4 STochFS Evaluation

In order to evaluate STochFS, we have selected feature selection algorithms that are accurate and highly characteristic. First, the LVF algorithm [11] is a Lagrange multiplier genetic algorithm based cardinality based feature selection algorithm. The Relief algorithm [9] is a distance-based feature selection algorithm based on the feature weights. We have also considered the decision algorithm, Focus [1] and ReliefD [8], the Focus algorithm has been effectively used in the data and ReliefD is the decision algorithm of Relief has a cardinality based data feature selection. Finally, the feature selection algorithm, we have added a distance based algorithm, aggregated, as described in section 3.2.

For each feature selection algorithm, we performed a series of experiments using the datasets: (C4.5, Naive Bayes and -Nearest Neighbors) and 13 datasets from the UCI Repository [3]: Credit (15 feature, 690 instances), Lab (16, 57), Vote (16, 435), Plasma Therapy (17, 339), Lymph (18, 148),

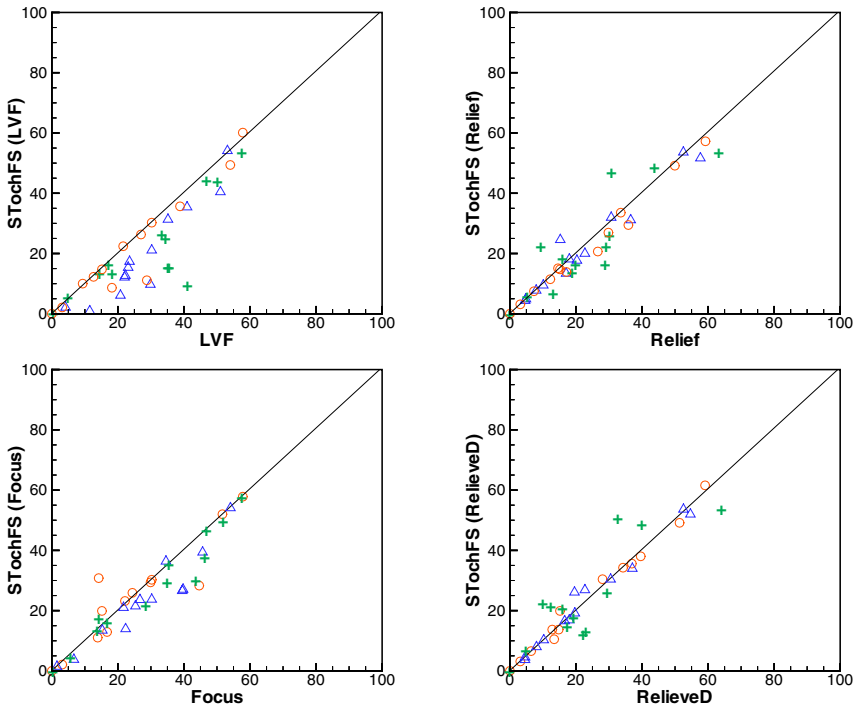


Fig. 4. Summary of the experimental results (error rates). Points under the line indicate that STochFS performed better than its underlying algorithm. Red circles indicate results for C4.5, blue triangles indicate results for Naive Bayes and green plus signs results for kNN.

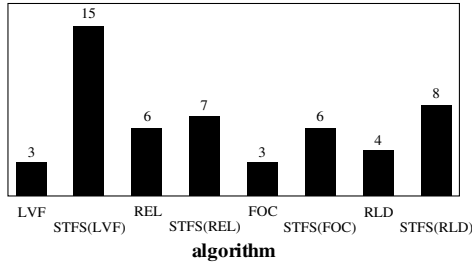


Fig. 5. Number of experiments (out of 39) which each algorithm performed the best or tied with the best. LVF = LVF, STFS(LVF) = STochFS using LVF, REL = Relief, STFS(REL) = STochFS using Relief, FOC = Focus, STFS(FOC) = STochFS using Focus, RLD = RelieveD and STFS(RLD) = STochFS using RelieveD.

M.h... (22, 8124), C.ic (23, 368), A... (25, 205), I...he,e (34, 351), S.bea... (35, 683), S.ice (60, 3190), S.a... (60, 208), A.d.i.g... (69, 226).

The following algorithms performed the best or tied with the best: LVF (3), STFS(LVF) (15), REL (6), STFS(REL) (7), FOC (3), STFS(FOC) (6), RLD (4), and STFS(RLD) (8). The STochFS variants using LVF, Relief, Focus, and RelieveD performed the best or tied with the best in 15, 7, 6, and 8 experiments, respectively. The STochFS variants using LVF, Relief, Focus, and RelieveD performed the best or tied with the best in 15, 7, 6, and 8 experiments, respectively. The STochFS variants using LVF, Relief, Focus, and RelieveD performed the best or tied with the best in 15, 7, 6, and 8 experiments, respectively.

The following algorithms performed the best or tied with the best: LVF (3), STFS(LVF) (15), REL (6), STFS(REL) (7), FOC (3), STFS(FOC) (6), RLD (4), and STFS(RLD) (8). The STochFS variants using LVF, Relief, Focus, and RelieveD performed the best or tied with the best in 15, 7, 6, and 8 experiments, respectively. The STochFS variants using LVF, Relief, Focus, and RelieveD performed the best or tied with the best in 15, 7, 6, and 8 experiments, respectively. The STochFS variants using LVF, Relief, Focus, and RelieveD performed the best or tied with the best in 15, 7, 6, and 8 experiments, respectively.

Both the STochFS variants using LVF and the STochFS variants using Relief performed the best or tied with the best in 15 and 7 experiments, respectively. The STochFS variants using Focus and RelieveD performed the best or tied with the best in 6 and 8 experiments, respectively. The STochFS variants using LVF, Relief, Focus, and RelieveD performed the best or tied with the best in 15, 7, 6, and 8 experiments, respectively.

³ To get to this number, we have tried different values for several datasets of small and medium sizes (up to 69 features) and the results showed that the STochFS performance is hurt, in several cases, if we use less than ten outcomes. Furthermore, using more than ten does not improve its performance in most situations.

Table 1. Score of the number of experiments (out of 39) each algorithm performed better within each significance level (calculated with the student's t-test). A score "A x B" for a certain algorithm f and significance level s means that STochFS performed better than f within s A times. Similarly, it also means that algorithm f outperformed STochFS B times within s .

	<0.001	<0.005	<0.01
STochFS vs LVF	19 x 0	4 x 0	4 x 1
STochFS vs Relief	12 x 3	3 x 0	4 x 2
STochFS vs Focus	7 x 3	6 x 1	3 x 0
STochFS vs RelieveD	6 x 6	2 x 1	2 x 2

generated the most significant features. Yet, Relief, a well-known characteristic selection method, performed better than STochFS in 4 out of 39 experiments. In addition, STochFS also aggregated the leading methods, outperforming Focus and RelieveD.

In order to evaluate the effectiveness of STochFS in selecting good feature subsets, we performed a blind experiment using the feature selection algorithms, where we identified each algorithm's best performance based on the accuracy of the aggregated features. In addition, we compared STochFS against LVF, Relief, Focus and RelieveD and the STochFS algorithm using the effectiveness of aggregated features. The results are shown in Figure 5, which shows the best performance of each algorithm. In addition, we identified the best performance of the aggregated features. Out of the 39 experiments (combination of 3 cases and 13 datasets), the accuracy of STochFS in selecting the best aggregated features was 36 cases (92.3%). This result can be compared to the best performance of each algorithm. The accuracy of the aggregated features was better than STochFS.

5 Conclusion

In this paper, we have evaluated the effectiveness of each feature selection algorithm in generating a set of combined features. The results show that STochFS, in addition to the best performance of Baggioli, has also performed better than each of the feature selection algorithms. The accuracy of the aggregated features was better than the accuracy of each algorithm. In addition, the accuracy of the aggregated features was better than the accuracy of each algorithm.

Therefore, we have shown that the effectiveness of the feature selection algorithms can be improved by using the aggregated features. In addition, STochFS achieved a better performance than each of the compared algorithms in selecting the best aggregated features.

The experimental results and discussed in this paper can guide each researcher in the design of the aggregated features. For example, each decision tree classifier can be used to evaluate the effectiveness of STochFS in

algorithm. In Proceedings of the 18th International Conference on Machine Learning, volume 1, pages 111–119, San Francisco, CA, 2001. Morgan Kaufmann.

References

1. H. Almuallim and T.G. Dietterich. Learning with many irrelevant features. In *Proceedings of the Ninth National Conference on Artificial Intelligence (AAAI'91)*, volume 2, pages 547–552, Anaheim, CA, 1991. AAAI Press.
2. J. Bala, K. DeJong, J. Huang, H. Vafaie, and H. Wechsler. Using learning to facilitate the evolution of features for recognizing visual concepts. *Evolutionary Computation*, 4(3):297–311, 1996.
3. C.L. Blake and C.J. Merz. UCI repository of machine learning databases, 1998. <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
4. L. Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
5. S. Das. Filters, wrappers and a boosting-based hybrid for feature selection. In *Proceedings of the Eighteenth International Conference on Machine Learning*, 2001.
6. Y. Freund and R.E. Schapire. Experiments with a new boosting algorithm. In *Proceedings of the Thirteenth International Conference on Machine Learning (ICML'96)*, pages 148–156, 1996.
7. M.A. Hall. Correlation-based feature selection for discrete and numeric class machine learning. In *Proceedings of the Seventeenth International Conference on Machine Learning*. Stanford University, CA, Morgan Kaufmann Publishers, 2000.
8. G.H. John, R. Kohavi, and K. Pfleger. Irrelevant features and the subset selection problem. In *Proceedings of the Eleventh International Conference on Machine Learning (ICML'94)*, pages 121–129, 1994.
9. K. Kira and L.A. Rendell. A practical approach to feature selection. In *Proceedings of the Ninth International Workshop on Machine Learning*, pages 249–256, Aberdeen, Scotland, 1992. Morgan-Kaufmann.
10. I. Kononenko, M. Robnik-Sikonia, and U. Pompe. *ReliefF for estimation and discretization of attributes in classification.*, pages 31–40. Artificial Intelligence: Methodology, Systems, Applications. IOS Press, 1996.
11. H. Liu and R. Setiono. A probabilistic approach to feature selection - a filter solution. In *Proceedings of the Thirteenth International Conference on Machine Learning (ICML'96)*, pages 319–327, 1996.
12. M. Richeldi and P. Lanzi. ADHOC: A tool for performing effective feature selection. In *Proceedings of the International Conference on Tools with Artificial Intelligence*, pages 102–105, 1996.
13. M. Sebban and R. Nock. A hybrid filter/wrapper approach of feature selection using information theory. *Pattern Recognition*, (35):835–846, 2002.
14. E.P. Xing, M.I. Jordan, and R.M. Karp. Feature selection for high-dimensional genomic microarray data. In *18th International Conference on Machine Learning*, pages 601–608, San Francisco, CA, 2001. Morgan Kaufmann.

Speeding Up Logistic Model Tree Induction

Marco Sumner^{1,2}, Eibe Fallett², and Mark Hall²

¹ Institute for Computer Science,
University of Freiburg,
Freiburg, Germany

sumner@informatik.uni-freiburg.de

² Department of Computer Science,
University of Waikato,
Hamilton, New Zealand

{eibe, mhall}@cs.waikato.ac.nz

Abstract. Logistic Model Trees have been shown to be very accurate and compact classifiers [8]. Their greatest disadvantage is the computational complexity of inducing the logistic regression models in the tree. We address this issue by using the AIC criterion [1] instead of cross-validation to prevent overfitting these models. In addition, a weight trimming heuristic is used which produces a significant speedup. We compare the training time and accuracy of the new induction process with the original one on various datasets and show that the training time often decreases while the classification accuracy diminishes only slightly.

1 Introduction

Logistic Model Tree (LMT) are based on the idea of combining classification and regression. The idea of logistic regression and decision trees has been shown to have LMT as a special case. The idea of combining classification and regression has been shown to be a useful technique [8]. However, the induction of LMT is a non-trivial task. This is due to the complexity of finding the logistic regression model. The Logistic Regression [6] problem is a well-known NP-complete problem. The decision boundary of a LMT is a linear combination of the logistic regression model. In this paper, we propose a heuristic method to induce LMT. The heuristic can be replaced by the AIC criterion [1] to improve accuracy. We also propose a weight trimming heuristic which has a significant speedup.

The rest of the paper is organized as follows. In Section 2 we give a brief overview of the original LMT induction algorithm. Section 3 describes the modified algorithm. Section 4 describes the algorithm. In Section 4, we also describe the modified algorithm and discuss the results. In Section 5 we describe the experimental results.

2 Logistic Model Tree Induction

The original LMT induction algorithm can be found in [8]. We give a brief overview of the algorithm. The algorithm is as follows. The algorithm is as follows.

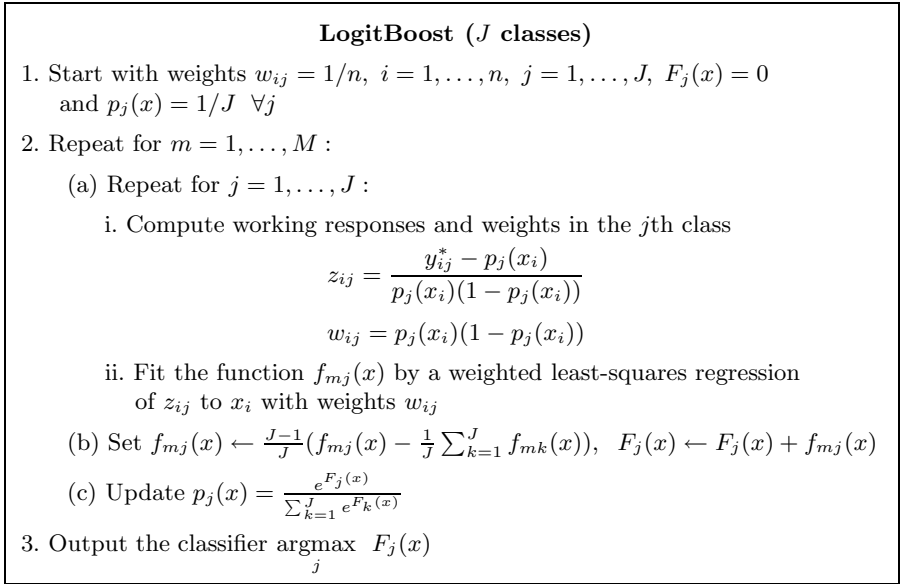


Fig. 1. LogitBoost algorithm

... We begin with each class having equal weight of the defining data and each class having equal weight in the LMT1 decision algorithm.

2.1 Logistic Regression

Let a logistic regression define the probability $Pr(G = j|X = x)$ of the J class label j for each x and define the probability $p_j(x)$ of the class label j in $[0, 1]$. The definition of the function

$$Pr(G = j|X = x) = \frac{e^{F_j(x)}}{\sum_{k=1}^J e^{F_k(x)}}, \tag{1}$$

where $F_j(x) = \beta_j^T \cdot x$. Note that the LMT1 algorithm has a similar form to the above definition of $p_j(x)$.

One choice of the LMT1 algorithm [6], shown in Figure 1, is to fit each class j to the data (x_i, y_{ij}) by a weighted least-squares regression of the data (x_i, y_{ij}) to the function $f_{mj}(x)$. Here, y_{ij}^* are the binary labels $\{0, 1\}$ which indicate the class label j for each x_i .

$$y_{ij}^* = \begin{cases} 1 & \text{if } y_i = j, \\ 0 & \text{if } y_i \neq j \end{cases}, \tag{2}$$

where y_i is the binary class label for x_i .

If each class f_{mj} is a linear function of x , then each f_{mj} is a linear function of x and the LMT1 algorithm is a linear combination of linear functions.

Automatic Iteration Termination. In the original design, the LMT heuristic based on the L_gB... (1) heuristic, failed to terminate because of a cycle in the heuristic. The heuristic based on the L_gB... (1) heuristic terminated because of a cycle in the heuristic.

A cycle in the heuristic based on the L_gB... (1) heuristic is a 1-2-3 cycle. The heuristic based on the L_gB... (1) heuristic terminated because of a cycle in the heuristic.

AIC is a heuristic based on the L_gB... (1) heuristic. The AIC is a heuristic based on the L_gB... (1) heuristic.

$$AIC = -\frac{2}{N} \loglik + 2 \frac{d}{N}, \tag{3}$$

where d is the number of parameters. If the number of parameters is large, the heuristic based on the L_gB... (1) heuristic is not as good as the heuristic based on the L_gB... (1) heuristic.

In this paper, we propose a heuristic based on the L_gB... (1) heuristic. The heuristic based on the L_gB... (1) heuristic is a heuristic based on the L_gB... (1) heuristic.

The heuristic based on the L_gB... (1) heuristic is a heuristic based on the L_gB... (1) heuristic. The heuristic based on the L_gB... (1) heuristic is a heuristic based on the L_gB... (1) heuristic.

Table 1. Training time and accuracy for SimpleLogistic and SimpleLogistic using weight trimming

Dataset	Training Time		Accuracy	
	SimpleLog.	SimpleLog. (WT)	SimpleLog.	SimpleLog. (WT)
vowel	77.94±23.59	39.67±12.72 ●	81.98±4.10	82.07±3.82
german-credit	7.97±1.94	6.79±1.55	75.37±3.53	75.35±3.48
segment	50.55±14.82	20.02±5.61 ●	95.10±1.46	86.71±25.67
splice	253.96±38.83	79.02±9.55 ●	95.86±1.17	95.87±1.09
kr-vs-kp	57.28±15.09	25.98±8.35 ●	97.06±0.98	97.07±0.92
hypothyroid	104.76±27.17	47.88±10.72 ●	96.61±0.71	96.55±0.72
sick	25.40±6.10	12.09±3.39 ●	96.68±0.71	96.63±0.70
spambase	119.28±18.73	43.19±4.38 ●	92.75±1.12	92.40±1.24
waveform	65.53±9.31	25.42±3.77 ●	86.96±1.58	86.90±1.55
optdigits	659.33±123.68	111.32±21.35 ●	97.12±0.67	97.17±0.67
pendigits	489.51±148.34	257.86±84.23 ●	95.44±0.62	95.51±0.61
nursery	266.51±25.56	119.19±11.36 ●	92.61±0.68	92.60±0.77
adult	2953.77±849.82	1866.15±344.05 ●	85.61±0.38	85.56±0.38

● statistically significant improvement

4.1 SimpleLogistic

We compare the accuracy of the SimpleLogistic algorithm using weight trimming (SimpleLogistic (WT)) against the accuracy of the SimpleLogistic algorithm using weight clipping (SimpleLogistic (WC)). The accuracy of SimpleLogistic (WT) is significantly higher than the accuracy of SimpleLogistic (WC) on 11 out of 16 datasets.

Weight Trimming in SimpleLogistic. From Table 1 it can be seen that SimpleLogistic (WT) is significantly more accurate than SimpleLogistic (WC) on 11 out of 16 datasets. The accuracy of SimpleLogistic (WT) is significantly higher than the accuracy of SimpleLogistic (WC) on the vowel, german-credit, segment, splice, kr-vs-kp, hypothyroid, sick, spambase, waveform, optdigits, pendigits, nursery, and adult datasets. The accuracy of SimpleLogistic (WT) is not significantly higher than the accuracy of SimpleLogistic (WC) on the vowel, german-credit, segment, splice, kr-vs-kp, hypothyroid, sick, spambase, waveform, optdigits, pendigits, nursery, and adult datasets.

FAM in SimpleLogistic. This section describes the accuracy of SimpleLogistic using FAM on 13 UCI datasets. The accuracy of SimpleLogistic using FAM is significantly higher than the accuracy of SimpleLogistic using FAM on 13 out of 13 datasets.

Table 2 shows the accuracy of SimpleLogistic using FAM on 13 UCI datasets. The accuracy of SimpleLogistic using FAM is significantly higher than the accuracy of SimpleLogistic using FAM on 13 out of 13 datasets. The accuracy of SimpleLogistic using FAM is significantly higher than the accuracy of SimpleLogistic using FAM on 13 out of 13 datasets.

4.2 Logistic Model Trees

We compare the accuracy of the Logistic Model Trees (LMT) algorithm against the accuracy of the SimpleLogistic algorithm using weight trimming (SimpleLogistic (WT)) on 13 UCI datasets. The accuracy of LMT is significantly higher than the accuracy of SimpleLogistic (WT) on 13 out of 13 datasets.

Table 2. Training time and accuracy for SimpleLogistic using cross-validation and FAM

Dataset	Training Time		Accuracy	
	SimpleLog. (CV)	SimpleLog. (FAM)	SimpleLog. (CV)	SimpleLog. (FAM)
vowel	77.94±23.59	6.87±0.31 ●	81.98±4.10	80.85±3.69
german-credit	7.97±1.94	0.59±0.05 ●	75.37±3.53	75.34±3.70
segment	50.55±14.82	3.42±0.45 ●	95.10±1.46	94.67±1.66
splice	253.96±38.83	77.48±3.69 ●	95.86±1.17	95.87±1.06
kr-vs-kp	57.28±15.09	6.69±0.37 ●	97.06±0.98	96.38±1.14 ○
hypothyroid	104.76±27.17	8.89±1.16 ●	96.61±0.71	95.89±0.65 ○
sick	25.40±6.10	1.57±0.14 ●	96.68±0.71	96.50±0.76
spambase	119.28±18.73	15.74±1.28 ●	92.75±1.12	92.69±1.19
waveform	65.53±9.31	7.75±0.39 ●	86.96±1.58	86.84±1.59
optdigits	659.33±123.68	135.61±26.47 ●	97.12±0.67	97.12±0.66
pendigits	489.51±148.34	59.43±1.58 ●	95.44±0.62	95.45±0.62
nursery	266.51±25.56	49.36±1.42 ●	92.61±0.68	92.58±0.68
adult	2953.77±849.82	381.92±10.13 ●	85.61±0.38	85.59±0.38

●, ○ statistically significant improvement or degradation

Table 3. Training time and accuracy for LMT and LMT using FAM and weight trimming

Dataset	Training Time		Accuracy	
	LMT	LMT (FAM+WT)	LMT	LMT (FAM+WT)
vowel	408.11±80.95	15.86±0.84 ●	94.06±2.40	93.56±2.94
german-credit	32.74±10.87	3.25±0.16 ●	75.50±3.65	71.83±3.40 ○
segment	143.75±52.64	10.58±1.77 ●	97.06±1.31	97.06±1.25
splice	785.51±202.14	71.55±1.42 ●	95.89±1.14	95.19±1.19 ○
kr-vs-kp	250.79±64.58	12.17±0.36 ●	99.64±0.33	99.57±0.37
hypothyroid	405.73±94.04	7.39±0.64 ●	99.54±0.36	99.61±0.30
sick	139.31±50.79	6.83±0.73 ●	98.95±0.58	98.93±0.62
spambase	746.71±123.57	54.93±1.65 ●	93.56±1.14	93.58±1.13
waveform	175.53±63.26	43.67±0.80 ●	86.86±1.60	86.49±1.52
optdigits	3162.37±781.49	133.15±7.08 ●	97.38±0.57	97.36±0.64
pendigits	3535.06±765.34	185.15±4.96 ●	98.58±0.33	98.73±0.33
nursery	634.96±85.82	72.44±7.08 ●	98.95±0.34	98.64±0.32 ○
adult	26935.85±9112.20	1429.93±54.76 ●	85.58±0.42	85.43±0.37

●, ○ statistically significant improvement or degradation

g ea e . . eed . . ec . ded . a 55 . . he h . . h . . id da a e . M . . c a a . eed . be . ee . 10 a d 25 . O he ge . a -c, edi , a ef . . , a d . . e . da a e . . a he . eed . a . . d 10 . . e . (10.1, 4.0, a d 8.8, e . ec i e) .

O . ge . a -c, edi , . . ice a d . . e . he . . di ed . e . i . ? c a i ca i . . e - f . . a ce . a ng i ca e ha . ha . f he . . ng i a . e . i . , a h . gh ge . a -c, edi . a he . ef . . a ce eb e ha e . ce . . O he . - i e . he . . di ed . e . i . . ef . . ed c . . e i i e . . i h he . . ng i a . e . i . .

A a c . . i g . . e e . e i e . . , e . . d i e . . c . . a e he . . di ed . e . i . . f . . g i c . . de . . ee . . b . . ed C4.5 dec i ee . F . . he c . . a i . . e ch . e AdaB . . [5] . i g 100 i e a i . . a d he LMT . e . i . . i g FAM a d eigh . . i . . i g . The . e . . ca . be . ee . 1 Tab e 4 . 100 i e a i . . a e a f . . a fe . . f he da a e . , b i g f . . 10 . . 100 i e a i e . . i a i e . e . i acc . ac i . . a ca e [8] .

The . ai . i g i . e . f he . . a g . i h . . i fa . . e . a , i ha . i gh ad a . age f . . he . . di ed LMT a g . i h . I . a fa . e . . 9 . f he 13 da a e . , i ha

Table 4. Training time and accuracy for AdaBoost using C4.5 with 100 iterations and LMT using FAM and weight trimming

Dataset	Training Time		Accuracy	
	AdaBoost	LMT (FAM+WT)	AdaBoost	LMT (FAM+WT)
vowel	29.32±0.38	15.86±0.84 ●	96.74±1.89	93.56±2.94 ○
german-credit	7.42±0.17	3.25±0.16 ●	74.40±3.23	71.83±3.40
segment	45.53±0.67	10.58±1.77 ●	98.58±0.76	97.06±1.25 ○
splice	11.89±5.46	71.55±1.42 ○	94.94±1.24	95.19±1.19
kr-vs-kp	21.14±6.44	12.17±0.36 ●	99.60±0.31	99.57±0.37
hypothyroid	19.07±11.46	7.39±0.64 ●	99.70±0.31	99.61±0.30
sick	49.40±2.32	6.83±0.73 ●	99.06±0.45	98.93±0.62
spambase	70.22±63.21	54.93±1.65	95.34±0.87	93.58±1.13 ○
waveform	463.38±4.18	43.67±0.80 ●	85.01±1.77	86.49±1.52 ●
optdigits	402.52±3.06	133.15±7.08 ●	98.55±0.50	97.36±0.64 ○
pendigits	274.59±2.72	185.15±4.96 ●	99.41±0.26	98.73±0.33 ○
nursery	24.90±0.48	72.44±7.08 ○	99.79±0.14	98.64±0.32 ○
adult	796.01±64.89	1429.93±54.76 ○	82.18±0.46	85.43±0.37 ●

●, ○ statistically significant improvement or degradation

...eed ... f ... a ... be ... ee ... 1 a d 2. The e ce ... a e he ic da a e (7.2) a d he a ef ... da a e (10.6). AdaB ... a fa e ... h ee da a e , he g ea e ... e e be i g ... he ... ice da a e (6.0). I e ... f c a i ca i ... acc , ac , he e LMT e ... e h b i ... e ... i a ... h e e ... ed i [8] f ... he ... i g a LMT a g ... h ... O ... i da a e. AdaB ... e f ... ed i g i ca ... be e , h e ... da a e. LMT a h e be e c a i e .

5 Conclusions

We ha e ... ed ... di ca i ... he Si ... eL g i c a g , i h e ... ed b LMT ha a e de i g e d ... i ... e , a i g i e . The ... e f AIC i ... e ad ... f c ... - a i d a i ... de e ... i e a a ... i a e ... be ... f L g i B ... i e a i ... e ... ed i a d a a i c ... eed ... I ... e ... ed i a ... a b ... i g i ca ... de c e a e i acc , ac i ... ca e ... he ... e f ... i g a d - a ... e g i c , e g e The ... e he ... i c f e i g h ... i ... i g c ... i e ... i ... ed he ... a i g i e h e ... a e c i g acc , ac a a .

The ... e f AIC a d e i g h ... i ... i g i LMT ha e e ... ed i ... a i g i e ... 55 i e fa e , ha he ... i g i a LMT a g ... h ... h i e , i ... ca e , ... i g i ca ... a e c i g c a i ca i ... acc , ac . The e e ... e e e a ... ed ... da a e ... f e a i e ... i e a d d i e ... i a i . We ... de ... ec he ... eed ... be e e g ea e ... a g e a d high-d i e ... i a da a e .

References

1. H. Akaike. Information theory and an extension of the maximum likelihood principle. In *Second Int Symposium on Information Theory*, pages 267–281, 1973.
2. C.L. Blake and C.J. Merz. UCI repository of machine learning databases, 1998. [www.ics.uci.edu/~mllearn/MLRepository.html].
3. L. Breiman, H. Friedman, J. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. Wadsworth, 1984.

4. Peter Bühlmann and Bin Yu. Boosting, model selection, lasso and nonnegative garrote. Technical Report 2005-127, Seminar for Statistics, ETH Zürich, 2005.
5. Yoav Freund and Robert E. Schapire. Experiments with a new boosting algorithm. In *Proc In. Conf on Machine Learning*, pages 148–156. Morgan Kaufmann, 1996.
6. Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Additive logistic regression: a statistical view of boosting. *The Annals of Statistics*, 38(2):337–374, 2000.
7. Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer-Verlag, 2001.
8. Niels Landwehr, Mark Hall, and Eibe Frank. Logistic model trees. *Machine Learning*, 59(1/2):161–205, 2005.
9. C. Nadeau and Yoshua Bengio. Inference for the generalization error. In *Advances in Neural Information Processing Systems 12*, pages 307–313. MIT Press, 1999.
10. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
11. I. H. Witten and E. Frank. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann, San Francisco, 2000.

A Random Method for Quantifying Changing Distributions in Data Streams

Haixun Wang¹ and Jian Pei²

¹ IBM T. J. Watson Research Center, USA
haixun@us.ibm.com

² Simon Fraser University, Canada
jpei@cs.sfu.ca

Abstract. In applications such as fraud and intrusion detection, it is of great interest to measure the evolving trends in the data. We consider the problem of quantifying changes between two datasets with class labels. Traditionally, changes are often measured by first estimating the probability distributions of the given data, and then computing the distance, for instance, the K-L divergence, between the estimated distributions. However, this approach is computationally infeasible for large, high dimensional datasets. The problem becomes more challenging in the streaming data environment, as the high speed makes it difficult for the learning process to keep up with the concept drifts in the data. To tackle this problem, we propose a method to quantify concept drifts using a universal model that incurs minimal learning cost. In addition, our model also provides the ability of performing classification.

1 Introduction

In this paper, we study *the distance between two data distributions* instead of two vectors or two sequences. Assume tuples in a training set D are drawn from an unknown distribution $F(\mathbf{x}, t)$. Each tuple is of the form (\mathbf{x}, t) , where \mathbf{x} is a vector and t is the class label of \mathbf{x} . The task of supervised learning or classification is to learn the unknown relationship between \mathbf{x} and t , that is, to find a model $f^*(\mathbf{x})$, such that the averaged difference between $f^*(\mathbf{x})$ and t is minimum.

We assume there are concept drifts in the unknown data distribution $F(\mathbf{x}, t)$. How do we quantify the concept drift by defining and computing the distance between the original dataset D and a new data set D' , which is drawn from the changed unknown distribution? Furthermore, how quantified changes can be used to tune the model $f^*(\mathbf{x})$ we learned before so that it maintains high accuracy on the changed data?

In the field of information theory, relative entropy, or the Kullback Leibler (K-L) divergence, has been suggested as an appropriate measure for comparing data distributions [5]. However, such methods are not computationally feasible for large, high dimensional datasets, or data coming from continuous streams. In the field of data mining, several works have studied how to *detect* changes of data distributions over streams and sequences [1,10]. However, more often than not, change detection only serves to trigger a costly learning process, and the change itself is not used to mend the current prediction model directly. Recently, several works [8,13] have studied how to update

the current model $f^*(\mathbf{x})$ in response to the concept drifts in data streams, for instance, by assimilating new instances in D' and forgetting old instances in D . These can be very costly undertakings since they do not handle changes directly on the probability distribution level, but rely on a lot of learning and re-learning.

We aim at devising an efficient method to measure distribution changes in high-dimensional, labeled datasets. We assign a *signature* to each dataset, and compare distribution changes by comparing the signatures. Furthermore, the signature should also enable us to make predictions.

2 A Model-Based Naive Approach

In this section, we introduce a naive but computationally feasible method for measuring distances between two datasets. We analyze the prediction error of this naive approach through bias/variance decomposition, and we study its impact on the distance measure. In the next section, we introduce a general approach based on the lessons learned here.

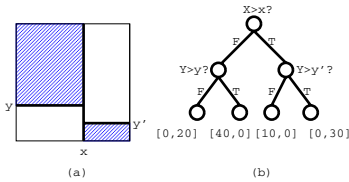


Fig. 1. Model-based description

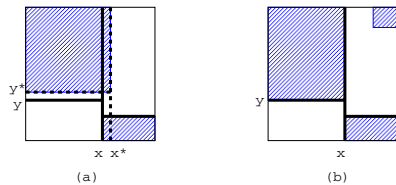


Fig. 2. Distribution Changes

2.1 Measuring Distribution Changes Using a Classification Model

Assume we are given a dataset D which consists of a set of tuples (\mathbf{x}, t) , where \mathbf{x} is a vector and t is the class label of \mathbf{x} . We learn a decision tree classifier T_D from D . The decision tree classifier T_D can be regarded as a summarization of the class distribution of dataset D . More specifically, let n_1, n_2, \dots, n_k be the leaf nodes of T_D . Each leaf node n_i is associated with a class distribution (number of objects belonging to each class). Together, (n_1, n_2, \dots, n_k) forms a special histogram of frequency counts.

For instance, in Figure 1(a) we show a two dimensional dataset where the shaded areas in the top-left and bottom-right corner are populated with objects of one class, and the rest of the area is populated with objects of the other class. In the rest of the paper, we assume the number of objects in an area is proportional to the size of the area.

From the dataset, we learn a decision tree classifier, which partitions the two dimensional space into 4 areas, each represented by a leaf node as shown in Figure 1(b). Each leaf node is associated with the number of objects of each class in that area. For instance, the second leaf node to the left represents the top-left area, where we assume $[40, 0]$ are the number of objects of the two classes in that area. All together, we can use the class distribution of the objects in the leaf nodes to describe the dataset. We call it the *signature* of the data:

$$([0, 20], [40, 0], [10, 0], [0, 30]) \tag{1}$$

Assume now there is some distribution change in the underlying dataset. In one case, the boundary of the shaded area moved from x to x^* horizontally and from y to y^* vertically, as shown in Figure 2(a).

We want to quantify the change using the model we learned from the original dataset. Here, we use the decision tree to classify the changed data set, and use the classification error to quantify the change. To a certain extent, the classification error represents the magnitude of the change, but certainly not the change itself. Because, for instance, datasets in Figure 2(a) and 2(b) will have the same classification error (compared with the original data set in Figure 1(a), they have the same amount of shaded area “out of the place”), but they have very different data distributions. Apparently, the error-based distance measure cannot be used to replace or tune the predictions made by the original decision tree for the changed data.

To ensure that the measure can represent, to a certain extent, the distribution of the change so that it can be used to help make predictions without learning a new model from the changed dataset, we simply ‘throw’ the objects in Figure 2(a) into the decision tree learned from the original dataset. The class distribution in the leaf nodes is now the signature of the changed dataset:

$$([0, 20], [38, 2], [10, 0], [2, 28]) \tag{2}$$

Now, the dataset in Figure 2(b) results in a different signature: $([0, 20], [40, 0], [10, 0], [4, 26])$, which means signatures are better than prediction errors in representing distributions.

Although we didn’t learn a decision tree from the new datasets, the signature, which combines the original decision tree structure and the new class distributions in the leaf nodes, give us some ability to make predictions. Take the dataset in Figure 2(a) and its signature Eq (2) for example. If a test object is classified into the 2nd leaf node to the left, the prediction that the object belong to the positive class will be the probability output $\frac{n_1}{n_1+n_2} = \frac{38}{38+2}$, where n_1 and n_2 are the number of positive and negative nodes in that leaf node respectively.

The signatures also enable us to quantify the differences between the two datasets. If we treat the signature as a vector, we can use any L_p metric to compute their distance. For example, the distance function Eq (3) between two signatures a and b is based on the Manhattan distance:

$$Dist_s(a, b) = \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^c \left| \frac{n_{a,j,k}}{N_a} - \frac{n_{b,j,k}}{N_b} \right| \tag{3}$$

where n is the number of leaf nodes, c is the number of different classes, $n_{a,j,k}$ is the number of nodes in the j -th leaf node that are of class label k , and N_a is the total number of objects in dataset a . For any two signatures a and b , we have $0 \leq Dist_s(a, b) \leq 1$.

This naive approach gives us the following benefits. First, it is computationally efficient to compare the differences of two data distributions. Second, the data descriptors can be used to make predictions. However, this naive method is also flawed.

2.2 Error Analysis

In the naive method, the model used to describe other datasets is partially learned from a dataset which may have a very different data distribution. This can result in significant prediction error and create problems for the distance measure. In this section, we first reveal such problems, then we use bias-variance decomposition to study their cause.

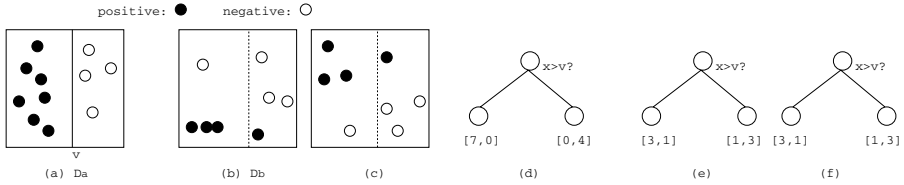


Fig. 3. A Greedy Learner

From D_a in Figure 3(a), we learn a decision tree, and we show the tree hierarchy in Figure 3(d). We then populate the leaf nodes of the decision tree with objects in other datasets. Figure 3(b) and 3(c) represent two very different data distributions. However, because of the tree structure learned from D_a , a same signature, $([3, 1], [1, 3])$, will be assigned to both datasets. Thus, the distance between the two very different distributions is 0. The signature is thus inaccurate because of the possible large variance introduced by training datasets such as D_a .

For similar reasons, using signatures assigned by the naive method for prediction is also flawed. We populate the leaf nodes of the decision tree learned from D_a in Figure 3(a) with objects in dataset D_b in Figure 3(b). This results in a signature of $([3, 1], [1, 3])$. Such a signature apparently has large prediction error – even when it is applied on D_b itself, the error can be as large as 25% under zero-one loss.

Clearly, this is due to the fact that D_a 's data distribution is very different from D_b 's. Decision trees are built in a divide-and-conquer, greedy manner, and in this case, there is no need to make a split on the Y axis for training set D_a , although such a split will result in the largest information gain as far as training set D_b is concerned. The difference of the two data distributions, combined with the greedy nature of the decision tree construction process, results in a large prediction error.

We observe samples (\mathbf{x}, t) drawn independently from some unknown distribution. We want to learn the unknown relationship between \mathbf{x} and t . That is, we want to find a function, $f^*(\mathbf{x})$, that minimizes a certain loss function $L(t, f^*(\mathbf{x}))$, where L can be zero-one loss, square loss, absolute loss, etc.

We use the notation $f^*(\mathbf{x}|D)$ to indicate that the prediction model f^* we learn depends on the training dataset D . We decompose the expected prediction error (EPE) into three terms: noise (σ^2), bias, and variance:

$$EPE(\mathbf{x}) = \sigma^2 + Bias(f^*(\mathbf{x}|D))^2 + Var(f^*(\mathbf{x}|D))$$

Let $E_D(f^*(\mathbf{x}|D))$ be the predicted value for sample \mathbf{x} averaged over all the training datasets. The variance can be expressed by:

$$E_D(E_D f^*(\mathbf{x}|D) - f^*(\mathbf{x}|D))^2$$

The variance term measures how sensitive the predicted value at \mathbf{x} is to random fluctuations in the training dataset. Traditionally, a model is learned from a training dataset D drawn from the data distribution we try to learn. In our case, we have two training sets, D_a and D_b . From D_a we learn the structure of the histogram (or equivalently the hierarchy of a decision tree), and from D_b we learn the data distributions within the structure or within the hierarchy. The variance can thus be expressed by:

$$E_{D_a, D_b}(E_{D_a, D_b} f^*(\mathbf{x}|D_a, D_b) - f^*(\mathbf{x}|D_a, D_b))^2$$

Since D_a might be drawn from a data distribution different from the distribution of D_b , which is the distribution we want to learn, by including both D_a and D_b in the condition, the variance is increased because of the added fluctuations.

3 A Universal Model

As discussed in the previous section, the majority of variance and bias is introduced due to training set D_a , from which we learn the structure of a histogram, or a hierarchy of a decision tree. Furthermore, it constitutes the major part of the learning cost. When the change of data distributions between D_a and D_b is non-trivial, the benefits of learning the tree structure from D_a becomes insignificant, since there is no guarantee that such a tree structure will fit the training dataset of D_b well. In this case, it becomes obvious that using an arbitrary tree structure not only serves the same purpose but at the same time eliminates the cost of learning such a structure.

Our goal is to find such an ‘arbitrary’ structure. It must be general and universal so that it can fit ‘any’ dataset D_b well, thus we can avoid the bias and variance component in the prediction error such as those introduced by one particular dataset D_a .

3.1 Distance by Random Signatures

A decision tree assigns a signature to a dataset. A signature can be regarded as a special histogram. Each bin, which corresponds to a leaf node in the decision tree, is ‘cut out’ or defined by the splitting conditions on the path from the root node to that leaf node. The learning procedure determines those conditions as well as their applying order through the computation of information gain.

Take the training set D_a in Figure 4 as an example. It is a two dimensional dataset with two class labels. From D_a , we learn a decision tree, which partitions the two dimensional space into a set of ‘bins’, each of which is in fact a leaf node in the decision tree. The signature is created by an *entropy-based partition*, since a decision tree is often constructed through the computation of information gain. Note that this learning procedure has super-linear complexity.

We propose to create signatures by randomly partitioning the multi-dimensional space into a set of bins. Figure 5 is such an example. The positions and the order of the splits are totally random, and instead of creating one histogram, we create multiple histograms, each of which is independently and randomly partitioned. In the following, we study two different ways of random partitioning.

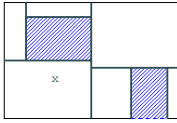


Fig. 4. A Decision Tree Histogram

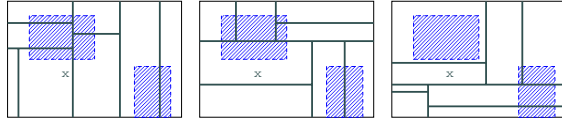


Fig. 5. Random Forest Histograms

Random Forest. We use the following procedure to create a random decision tree for a training dataset D .

1. $partition(D)$: randomly pick an unused attribute to partition D into D_1, \dots, D_n ;
2. for each partition D_i ($1 \leq i \leq n$), recursively invoke $partition(D_i)$ till the k -th recursive level.

We repeat this process N times to create a forest of N random trees [3]. Each tree defines a signature, and the random forest consists of N signatures for the dataset.

Random Histograms. We use the following procedure to create a random histogram for a training dataset D .

1. Randomly pick k attributes, a_1, \dots, a_k , as well as one value for each attribute, such that $\{a_1 = v_1, \dots, a_k = v_k\}$ defines a bin in the histogram.
2. Repeat the above step M times so that we have a histogram of M bins.

We repeat the above process N times to create N random histograms.

Each of the above methods creates N random structures. Given a dataset D_x , we populate the random structures with objects in D_x , which results in N signatures $S_{x,1}, \dots, S_{x,N}$ for D_x . We use the same random structures for all datasets. Clearly, for any two datasets, D_a and D_b , signatures $S_{a,i}$ and $S_{b,i}$ have the same number of bins and each bin defines the same subspace in the multi-dimensional data space. We then define the distance between two datasets D_a and D_b as: $Dist(D_a, D_b) = \frac{1}{N} \sum_{i=1}^N Dist_s(S_{a,i}, S_{b,i})$, where $Dist_s$ is the distance between two signatures defined in Eq (3), and we have $0 \leq Dist(D_a, D_b) \leq 1$.

The difference between this method and the naive method is that in this method, i) the structure of a signature does not rely on one dataset (which is known as D_a in the naive method), and ii) instead of having one signature, it uses multiple signatures. As will be discussed in detail in the following sections, the multiple random signatures is capable of ‘fitting’ any dataset, which means the distance metric and the prediction model will have high accuracy.

3.2 Classification by Random Signatures

A signature is composed of a set of histograms, each of which can be expressed by a vector $[n_1, \dots, n_c]$, where n_i is the number of objects that belong to class i .

The signature is used for prediction: an testing object that falls into a bin with class histogram $[n_1, \dots, n_c]$ is classified to be of class i if $i = \text{arg} \max_i \frac{n_i}{\sum n_i}$. However, a random signature is often a “weak” classifier.

The weakness of a single random signature can often be averted as our random methods create N signatures for a training dataset. The final prediction is a voted combination of all signatures. In other words, each signature is a classifier, and the N signatures form a classifier ensemble.

Combining an ensemble of classifiers is an established research area [2,6,12]. Particularly, for random forests, the prediction accuracy is shown to be no less than that of normal decision trees. Although each random signature is possibly a very “weak” classifier, it has been shown that if each classifier in the ensemble is independent in the production of its error, the expected error of the ensemble can be reduced to zero as the number of the classifiers goes to infinity [7].

3.3 Signatures’ Structural Diversity

Whether the signature-based distance metric and prediction model are meaningful depends on whether the random signatures can “fit” any dataset. The strength of an ensemble comes from its diversity [9]. In this section, we discuss how to guarantee signatures’ structural diversity.

In an ensemble, a classifier is valuable if it disagrees on some inputs with the other classifiers. Building a diverse ensemble in which each hypothesis is as different as possible is important to an ensemble method. Normally, diversity is measured by prediction disagreements among ensemble classifiers. In our case, random structures are created without a training dataset, which means we can only measure diversity by directly studying the differences of their internal structures. In a signature, each bin corresponds to a set of attribute values. We use the number of different attribute combinations as a measure of diversity. Let A be the number of attributes of the datasets. For simplicity, in our discussion we assume each attribute has v unique values.

- In a *random forest*, each tree of height k has v^k leaf nodes. The path from the root node to any leaf node has $k - 1$ edges. Thus, the diversity of attribute combinations in one random tree is at most $\min(v^{k-1}, \binom{k-1}{A})$. In the worst case, all leaf nodes (bins) share one attribute combination. Furthermore, attribute combinations may be correlated.
- For *random histograms*, each bin is defined independently by k attribute values. To compare with the above methods, we create v^k bins. The diversity can be as high as $\min(v^k, \binom{k}{A})$. In the worst case, all bins share one attribute combination. This occurs when all attributes are used ($k = A$), or each random selection returns the same set of attributes.

In summary, random histograms provide the most diverse set of attribute combinations with low correlation.

Our second question is how many bins should we keep in each random structure? We answer this question for random histograms. For random histograms, the number of attribute combinations is at most $\min(v^k, \binom{k}{A})$. Note that $\binom{k}{A}$ reaches maximum

when $k = A/2$. Thus, when $v^{A/2} > \binom{A/2}{A}$, we shall use $k = A/2$ attributes for random histograms; otherwise, we shall use k attributes where k satisfies $v^k \geq \binom{k}{A}$ and $v^{k-1} < \binom{k-1}{A}$.

4 Conclusion

The ability to quantify the similarity between two datasets is important to many applications, especially data stream applications that deal with time-changing data distributions. Statistical methods, such as K-L divergence and Kriging, are usually not computationally feasible for large, high speed datasets. In this paper, we propose a new approach based on the theory of random forests and classifier ensemble. To measure the difference between two data distributions, our approach measures the difference between the models derived from the datasets. To do this, we must use models that can truthfully represent the dataset, and models that can be trained efficiently. The models we propose for this purpose is the random histograms. The random histograms assign datasets signatures, which serve for two purposes: i) to measure distance between datasets by directly comparing signatures; and ii) to perform classification.

References

1. Charu C. Aggarwal. A framework for diagnosing changes in evolving data streams. In *SIGMOD*, 2003.
2. Eric Bauer and Ron Kohavi. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine Learning*, 36(1-2):105–139, 1999.
3. L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
4. Sergey Brin. Near neighbor search in large metric spaces. In *VLDB*, Switzerland, 1995.
5. Thomas M. Cover. *Elements of Information Theory*. Wiley-Interscience, 1991.
6. Yoav Freund and Robert E. Schapire. Experiments with a new boosting algorithm. In *ICML*, pages 148–156, 1996.
7. L. Hansen and P. Salamon. Neural network ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12:993–1001, 1990.
8. G. Hulten, L. Spencer, and P. Domingos. Mining time-changing data streams. In *SIGKDD*, pages 97–106, San Francisco, CA, 2001. ACM Press.
9. A. Krogh and J. Vedelsby. Neural network ensembles, cross validation, and active learning. In *Advances in Neural Information Processing Systems*, volume 7, pages 231–238. MIT Press, 1995.
10. Junshui Ma and Simon Perkins. Online novelty detection on temporal sequences. In *SIGKDD*, 2003.
11. M. A. Oliver and R. Webster. Kriging: a method of interpolation for geographical information systems. *International Journal Geographic Information Systems*, 4(3), 1990.
12. Kagan Tumer and Joydeep Ghosh. Error correlation and error reduction in ensemble classifiers. *Connection Science*, 8(3-4):385–403, 1996.
13. Haixun Wang, Wei Fan, Philip S. Yu, and Jiawei Han. Mining concept-drifting data streams using ensemble classifiers. In *SIGKDD*, 2003.

base, de la g. e b die . CLIQUE, DOC a d SUBCLU ee de . e c . e . 1 e e . b ace [4,5,6]. H e e , he d . . . ide QFI c . e . i g b . h . . e . ic a d ca eg . i ca a . i b e a d c . e . . . di g . a 1 - a a e a d h e . - ec a g a c . e . 1 a e cie . a . e . The a . . . a che . . . i e a - 1 a i e a . . cia 1 . . e ha e adde . ed hi 1 . e [7,8]. H e e , he di c e i e each a . i b e i h . c . i d e i g he de e de c a . . g a . i b e , a d h . ca . e . . i . he af . e e . i . ed f ag e a i . . .

I h i . d , a . . e a d e cie . a . . . a ch . . . he e ha . i e , a 1 - a a e a d h e . - ec a g a . b ace c . e . i g 1 . . . ed . M . e . e , he c . b i ed . e . f h i c . e . i g a d CAEP i e a a ed . The . . . ed c . e . i g a g . i h ha a e e i e . . c . e . Whie h i 1 . i a . . SUBCLU , . . . a . . . a ch ca de i e c . e . . . b . h . . e . ic a d ca eg . i ca i e . . , a d he . . . e . ic i e . ha i g 1 e , a . a e ca . be . . . ce . ed .

2 CAEP

CAEP i b i e e . a i ed a . . . [3]. The . a i i g ha e . f CAEP c . e . i . . . f . . . ce . e . The . . . i . de i e a . . e b die . Le he f a i e . e a b D_{cl} be $support_{D_{cl}}(a) = |\{t \in D_{cl} | a \in t\}| / |D_{cl}|$. F . e . e . cl , a . e . f QFI , $LQFI(cl)$, i h i c h e e . i e . e a a i e $support_{D_{cl}}(a) \geq min.sup$, i de i ed f . . D_{cl} . S b e e . , f . e . e . $a \in LQFI(cl)$, he f . . i g f . a ca . cl i ca c a ed . Le $D_{cl} = D - D_{cl}$ be he . . . e . i . a ce . f cl .

Growth rate:

If $support_{\bar{D}_{cl}}(a) \neq 0$, $growth_rate_{\bar{D}_{cl} \rightarrow D_{cl}}(a) = \frac{support_{D_{cl}}(a)}{support_{\bar{D}_{cl}}(a)}$,

If $support_{\bar{D}_{cl}}(a) = 0$ a d $support_{D_{cl}}(a) \neq 0$, $growth_rate_{\bar{D}_{cl} \rightarrow D_{cl}}(a) = \infty$,

O he . i e $growth_rate_{\bar{D}_{cl} \rightarrow D_{cl}}(a) = 0$.

Whe he g . h a e . f a i . . . e ha a $\rho (> 1)$, . . . , $growth_rate_{\bar{D}_{cl} \rightarrow D_{cl}}(a) \geq \rho$, a i ca ed a (EP) a d e ec ed a a . e b d . he e i . head ha he ca . cl , . . . , $a \Rightarrow cl$. Le $LEP(cl)$ be a . e . f a EP . e ec ed f . . $LQFI(cl)$. de . h i . ea . e . The . de . i g . i c i e he e i . . . e ec . he . e b die ha i g he . e g h . di e e . i a e he ca . cl f . . he . he . . E e . if he . e c . . de ce i high i D_{cl} , he . e ca . a ch . . a . i . a ce i D_{cl} . S ch . e a e ea f . ca i ca i . . .

The ec . d . ce . i . de i e a F i . , he . e g h f a EP a ba ed . . he . e a i e di e e ce be ee $support_{D_{cl}}(a)$ a d $support_{\bar{D}_{cl}}(a)$ i . i . . d ced a $support_{D_{cl}}(a) / (support_{D_{cl}}(a) + support_{\bar{D}_{cl}}(a)) = growth_rate_{\bar{D}_{cl} \rightarrow D_{cl}}(a) / (growth_rate_{\bar{D}_{cl} \rightarrow D_{cl}}(a) + 1)$. The f . . i g e . e e . he . . i b i i . f t . . be ca i ed i . cl b EP i $LEP(cl)$.

Aggregate score:

$$score(t, cl) = \sum_{a \subseteq t, a \in LEP(cl)} \frac{growth_rate(a)}{growth_rate(a) + 1} * support_{D_{cl}}(a). \tag{1}$$

Beca e he . . be . f EP f . each cl . a . . . be ba a ced, i . a ce . a ge h i g h e . c . e f . . . e ca . e . A ba e c . e i i . . d ced . e i i a e h i b i a .

Base score:

$base_score(cl)$ is the edialfa aggregation of $\{score(t, cl) | t \in D_{cl}\}$. The edge has the $base_score(cl)$, $growth_rate(a)$ and $support_{D_{cl}}(a)$ based on the algorithm. Given a set, the aggregation of cl , $score(t, cl)$, is calculated from the edge and E.(1). The, is based on $base_score(cl)$ and the edge is defined by a formula.

Normalized score:

$$norm_score(t, cl) = \frac{score(t, cl)}{base_score(cl)}$$

cl has the algorithm based on the edge and the call of t . The edge has the $LQFI(cl)$ of cl , the call of the edge of CAEP $O(N)$ where $N = |D|$, is calculated by the algorithm.

3 Mining Rule Bodies of CARs

3.1 Levelwise Subspace Clustering

Figure 1 shows the clustering process. The edge is defined by the edge and the edge. The edge is defined by the edge and the edge. The edge is defined by the edge and the edge. The edge is defined by the edge and the edge.

Let t and t' be the edge has the edge and the edge p in the edge q and q' is the edge. The Δ_p is defined by $N_{\Delta_p}(t) = \{t' \in D_{cl} | Dist_p(q, q') \leq \Delta_p\}$ where Δ_p is the edge. If the edge q and q' is the edge, the $Dist_p(q, q') = 0$, the edge $Dist_p(q, q')$ is the distance between the edge. The edge $t \in D_{cl}$ is called a p if $N_{\Delta_p}(t)$ contains at least $MinPts$ edge, $|N_{\Delta_p}(t)| \geq MinPts$. When a set of edge t is called $N_{\Delta_p}(t')$ for the edge t' , t and t' have a

Definition 1 (Density-Connected Set). $C \subseteq D_{cl}$ is a density-connected set if for every $p \in C$, $N_{\Delta_p}(p) \cap C \neq \emptyset$.

Definition 2 (Dense Cluster). $C^S \subseteq D_{cl}$ is a dense cluster if C^S is a density-connected set and $S \subseteq C^S$ for every $p \in S$, $D_{cl} \setminus C^S$ is not a density-connected set.

Definition 3 (Quantitative Frequent Itemset). $C^S \subseteq D_{cl}$ is a quantitative frequent itemset if $S \subseteq C^S$ and $a(C^S) = \{ \langle p : q \rangle \mid p \in S, q = [_1_p(C^S), _a_p(C^S)] \}$ where $_1_p(C^S) = _a_p(C^S)$ and $_a_p(C^S) = _a_p(C^S)$. $support_{D_{cl}}(a(C^S)) \geq minsup$ and $|C^S| \geq minsup$.

A QFI is defined as a set of edge and the edge. The edge has the edge and the edge. The edge has the edge and the edge. The edge has the edge and the edge.

Table 1. An example of transaction data set of $cl = \text{Houseowner}$; $D_{\text{Houseowner}}$

$t_1 = (\{ \langle \text{Age} : [20, 23] \rangle, \langle \text{Child} : [2, 3] \rangle, \langle \text{NumCars} : [2, 2] \rangle, \text{Houseowner})$
$t_2 = (\{ \langle \text{Age} : [30, 30] \rangle, \langle \text{Child} : [4, 5] \rangle, \langle \text{NumCars} : [1, 1] \rangle, \langle \text{Savings} : [10K, 10K] \rangle, \text{Houseowner})$
$t_3 = (\{ \langle \text{Age} : [30, 30] \rangle, \langle \text{Child} : [2, 2] \rangle, \langle \text{NumCars} : [5, 5] \rangle, \langle \text{Savings} : [11K, 11K] \rangle, \text{Houseowner})$
$t_4 = (\{ \langle \text{Age} : [30, 35] \rangle, \langle \text{Child} : [5, 5] \rangle, \langle \text{NumCars} : [1, 1] \rangle, \text{Houseowner})$
$t_5 = (\{ \langle \text{Age} : [35, 37] \rangle, \langle \text{Child} : [2, 2] \rangle, \langle \text{NumCars} : [2, 2] \rangle, \langle \text{Savings} : [5K, 5K] \rangle, \text{Houseowner})$
$t_6 = (\{ \langle \text{Age} : [36, 39] \rangle, \langle \text{Child} : [2, 2] \rangle, \langle \text{NumCars} : [2, 3] \rangle, \text{Houseowner})$

Table 2. Process of levelwise subspace clustering of $D_{\text{Houseowner}}$

1-QFIs $(\{ \langle \text{Age} : [30, 39] \rangle, \{t_2, t_3, t_4, t_5, t_6\} \}, (\{ \langle \text{Child} : [2, 5] \rangle, \{t_1, t_2, t_3, t_4, t_5, t_6\} \})$ $(\{ \langle \text{NumCars} : [1, 3] \rangle, \{t_1, t_2, t_4, t_5, t_6\} \}, (\{ \langle \text{Savings} : [10K, 11K] \rangle, \{t_2, t_3\} \})$
2-QFIs $(\{ \langle \text{Age} : [30, 39] \rangle, \langle \text{Child} : [2, 2] \rangle, \{t_3, t_5, t_6\} \})$ $(\{ \langle \text{Age} : [30, 35] \rangle, \langle \text{Child} : [4, 5] \rangle, \{t_2, t_4\} \})$ $(\{ \langle \text{Age} : [30, 39] \rangle, \langle \text{NumCars} : [1, 3] \rangle, \{t_2, t_4, t_5, t_6\} \})$ $(\{ \langle \text{Age} : [30, 30] \rangle, \langle \text{Savings} : [10K, 11K] \rangle, \{t_2, t_3\} \})$ $(\{ \langle \text{Child} : [2, 5] \rangle, \langle \text{NumCars} : [1, 3] \rangle, \{t_1, t_2, t_4, t_5, t_6\} \})$
3-QFIs $(\{ \langle \text{Age} : [35, 39] \rangle, \langle \text{Child} : [2, 2] \rangle, \langle \text{NumCars} : [2, 3] \rangle, \{t_5, t_6\} \})$ $(\{ \langle \text{Age} : [30, 35] \rangle, \langle \text{Child} : [4, 5] \rangle, \langle \text{NumCars} : [1, 1] \rangle, \{t_2, t_4\} \})$

Lemma 1 (Monotonicity). $\forall T \subseteq S, a(C^S) \subseteq a(C^T) \implies a(C^S) \subseteq a(C^T)$

Because C^S is a derived closed set, $\forall p \in S, p$ is a derived closed set. $\forall p \in T, a$ and hence $C^S \subseteq C^T$. Therefore, $\forall p \in T, [p]_p(C^S), [p]_p(C^S) \subseteq [p]_p(C^T), [p]_p(C^T)$, and $a(C^T)$ is a derived set of $a(C^S)$. ■

Accordingly, we can reach a local maximum for each a QFI. We can see if there are any better results in Table 1. Each transaction t_i contains a set of attributes $cl = \text{Houseowner}$. We find the characteristic of $\Delta_{\text{Age}} = 5, \Delta_{\text{Child}} = 1, \Delta_{\text{NumCars}} = 1, \Delta_{\text{Savings}} = 1K, \text{MinPts} = 1$ and $\text{minsup} = 2$. First, we find each t_i are a regular itemset, and the itemset $\{ \langle \text{Age} : [30, 39] \rangle, \{t_2, t_3, t_4, t_5, t_6\} \}$ has a $\text{minsup} = 5$ and $\Delta_{\text{Age}} = 5$, and $\text{minsup} = 5$ and minsup . This 1-QFI is a local maximum. (TID-L1) is a local maximum in Table 2. Each itemset has a 1-QFI in the table.

If there are k itemsets, we can find each k -QFI ($k > 1$) and k -QFI. We can find $TID - List$ and $TID - List'$ and $TID - List^c$. A k -QFI is a $(k-1)$ -QFI and k -QFI.

Definition 4 (Candidate-Generation).

Join Phase $(k-1)$ \dots $k-2 \dots$
 $((k-1) - QFI = \{ \langle p_1 : q_1 \rangle, \langle p_2 : q_2 \rangle, \dots, \langle p_{k-2} : q_{k-2} \rangle, \langle p_{k-1} : q_{k-1} \rangle, TID - List \},$
 $((k-1) - QFI' = \{ \langle p_1 : q'_1 \rangle, \langle p_2 : q'_2 \rangle, \dots, \langle p_{k-2} : q'_{k-2} \rangle, \langle p_k : q'_k \rangle, TID - List' \},$
 $(candidate - k - QFI = \{ \langle p_1 : q^c \rangle, \dots, \langle p_{k-1} : q^c_{k-1} \rangle, \langle p_k : q^c_k \rangle, TID - List^c \}.$

<p>QFI-Count(<i>candidate</i> - <i>k</i> - QFI, <i>TID</i> - <i>List</i>^c);</p> <p>(1) $k - QFIS = \phi, TIDLS = \phi;$</p> <p>(2) If $TID - List^c < minsup$ return $k - QFIS;$</p> <p>(3) $S = \{p \mid \langle p : q \rangle \in candidate - k - QFI, p \text{ is numeric.}\};$</p> <p>(4) $TIDLS.temp = \{TID - List^c\};$</p> <p>(5) while $TIDLS \neq TIDLS.temp$ do begin</p> <p>(6) $TIDLS = TIDLS.temp;$</p> <p>(7) forall $p \in S$ do begin</p> <p>(8) $TIDLS.temp = MDCS(TIDLS.temp, p);$</p> <p>(9) end</p> <p>(10) end</p> <p>(11) forall $TID - List \in TIDLS$ do begin</p> <p>(12) $k - QFIS = k - QFIS + (QFI(S, TID - List), TID - List);$</p> <p>(13) end</p> <p>(14) return $k - QFIS;$</p>	<p>(1) For each numeric attribute, create an index list sorted with the ascending order of D. Sort items in each $t \in D$ lexicographically.</p> <p>(2) $L_1 = \{(1 - QFI, TID - List)\};$</p> <p>(3) for ($k=2; L_{k-1} \neq \phi; k++$) do begin</p> <p>(4) $C_k = \{(candidate - k - QFI, TID - List^c)\} = Extended - Candidate - Generation(L_{k-1});$</p> <p>(5) forall ($candidate - k - QFI, TID - List^c \in C_k$) do begin</p> <p>(6) $L_k = L_k \cup \{QFI - Count(candidate - k - QFI, TID - List^c)\}$</p> <p>(7) end</p> <p>(8) end</p> <p>(9) Answer $L = \bigcup_k L_k;$</p>
--	---

Fig. 1. Algorithm of QFI-Count

Fig. 2. Entire algorithm

$q_i^c \dots q_i \cap q_i' \dots i = 1, \dots, k - 2,$
 $q_{k-1}^c = q_{k-1}, q_k^c = q_k' \dots TID - List^c = TID - List \cap TID - List'$
 $q_i^c = \phi, \dots (k - 1)$

Prune Phase $(k - 1) \dots s \dots k \dots (k - 1)$

$$\forall \langle p_i : q_i^c \rangle \in s, \exists \langle p_i : q_i \rangle \in (k - 1) - QFI, q_i^c \cap q_i \neq \phi, \quad (2)$$

$k \dots TID - List^c \dots C^S \dots |S| = k \dots k$

This phase has been established. Lemma 1. A fact that $q_i^c \cap q_i \neq \phi$ in (2), the attribute has satisfied $(k - 1) - QFI$ has been the $minsup$ attribute. It is eligible. The candidate $k - QFI$ is added to the candidate list. In Table 2, a candidate $e - 2 - QFI, \{\langle Age : [30, 39] \rangle, \langle Child : [2, 5] \rangle\}$ in $TID - List^c = \{t_2, t_3, t_4, t_5, t_6\}$ is defined for $1 - QFI, \{\langle Age : [30, 39] \rangle\}$ and $\{\langle Child : [2, 5] \rangle\}$. This attribute has been

... Fig.1 defined for $C^S = TID - List$ and the candidate $k - QFI$, if the attribute p has been defined for C^S based on Definition 2 and 3. In the definition of (7) - (9), a candidate C is checked for p in C^S and a fact that $MDCS$ is generated for Δ_p and $MinPts$. Moreover C can be found for the candidate C^S . $MDCS$ is added to C and p is added to C and added to $TIDLS.temp$ and the candidate C has a gain of $minsup$ in $MDCS$. This data is added to the candidate list (5) - (10), for each C in C^S the candidate C^S is defined for the candidate C and the attribute p is added to C . In the definition (11) - (13), each QFI candidate $C^S = TID - List$ is added to QFI and added to the candidate list. In the candidate $e - 2 - QFI, \{\langle Age : [30, 39] \rangle, \langle Child : [2, 5] \rangle\}$

1 h $TID - List^c = \{t_2, t_3, t_4, t_5, t_6\}$ 1 g i e . . . h i QFI-C . . . I he 1 ide . . . , $MDCS$ de 1 e $TIDLS.temp = \{\{t_2, t_3, t_4, t_5, t_6\}\}$. . . Age. Ne , 1 de 1 e $TIDLS.temp = \{\{t_3, t_5, t_6\}, \{t_2, t_4\}\}$. . . Child. F he a ica 1 . . . f $MDCS$ d . . . cha ge $TIDLS.temp$. Si ce he 1 e f ca dida e a e . . . e a . $minsup = 2$, . . . 2-QFI , ($\langle Age : [30, 39] \rangle$, $\langle Child : [2, 2] \rangle$), $\{t_3, t_5, t_6\}$ a d ($\langle Age : [30, 35] \rangle$, $\langle Child : [4, 5] \rangle$), $\{t_2, t_4\}$), a e de 1 ed.

3.2 Deriving QFIs of Numeric and Categorical Items

Ca dida e-Ge e a 1 . 1 e e ded . de 1 e QFI c . . 1 1 g f . . e ica d ca - eg , ica 1 e . . . The ca eg , ica 1 e . . 1 he 1 ed 1 e e a e g i e 1 he a e a a 1 he A , 1 Tid a g , 1 h . I he 1 ha e f De . 1 1 . 4, if $q_i^c = \phi$ f . . . e . . e ica 1 e . . . $q_i \neq q'_i$ f . . . e ca eg , ica 1 e , he . . . g i e (k-1)-QFI a e . . . 1 ed. O he 1 e he a e 1 ed a $q_i^c = q_i \cap q'_i$ f , each . . e ica 1 e 1 e a d $q_i^c = q_i = q'_i$ f , each ca eg , ica 1 e . I he . . e ha e, he c . d i . . $q_i^c = q_i$ f , a ca eg , ica 1 e 1 a 1 ed 1 add 1 . . . $q_i^c \cap q_i \neq \phi$ f , a . . e ica 1 e 1 E .(2). The a g , 1 h QFI-C . . . f Fig.1 1 a . a e ed. Whe he ca dida e-k-QFI c . . 1 . . f ca eg , ica 1 e . . . , he . . . f . . (5) . (10) 1 . . 1 ed, a d $TIDLS = TIDLS.temp$ 1 a 1 ed. The f c 1 . QFI a . e (12) 1 a . a e ed. F a ca eg , ica a 1 b e p_i , 1 . a e 1 e . . be $q_i^c = q_i = q'_i$.

The e 1 e a g , 1 h . . de 1 e QFI f . . D 1 1 dica ed 1 Fig.2. Re 1 ed a a e e . a e Δ_p f , a . . e ica 1 b e , $MinPts$ a d $minsup$. F 1 . , . . e 1 de 1 . a e c e a ed f , he e c i e . . ce 1 g 1 E e ded-Ca dida e-Ge e a 1 . a d QFI-C . . . S be e e , a QFI a e c . . ed 1 L b he ada a 1 . f he A , 1 Tid A g , 1 h . I he 1 e e a 1 , he 1 e . ed 1 de 1 g ($t_i, \{candidate - k - QFI\}$) f . . each t_i . 1 . c . a 1 g ca dida e-k-QFI 1 . ed 1 . ead f ($candidate - k - QFI, TID - List^c$) . 1 1 a . . he . a da d A , 1 Tid. Thi a . . ach 1 a 1 ed . D_{cl} f e e . ca . cl . de 1 e $LQFI(cl)$, e 1 ed b CAEP de c 1 b ed 1 he . e 1 . . ec 1 . .

4 Experimental Evaluation

4.1 Computational Efficiency

The . . . e e 1 e a 1 he de 1 a 1 . . f QFI f , CAR' b die . Th . , 1 . c . . a 1 . a e c i e c 1 e a a ed b . 1 g a 1 cia da a e . F 1 . , a e f f . eed 1 e . . , SSI , 1 . a d . . ge e a ed he e $r_n\%$ f he a e . . e ica d he . e ca eg , ica . Sec d, a e f f eed QFI , $SQFI$, 1 ge e a ed b . a d . . eec 1 g eed 1 e . f . . SSI . The 1 e f each QFI 1 de e . 1 ed b . . f . . a d . d i . b 1 . ha 1 g 1 . a e age a $\overline{|QFI|}$. Th d, a e f 1 . a ce (. a . - ac 1 . .) D 1 ge e a ed he e each 1 . a ce t 1 . ade b . a d . . eec 1 g a QFI f . . $SQFI$ a d f , he . a d . . add 1 ge . a $2\overline{|QFI|}$. eed 1 e . a e f . . SSI 1 he a e age. F 1 a . , he a e f . . e ica 1 e . . 1 each t a e d i . . ed b 1 . . d c 1 g Ga . . ia . . 1 e ha 1 g 5% a . 1 de . O , a g , 1 h 1 e e ed . a Pe . 1 . 4.2.7 GH PC 1 h 2GB RAM. The defa . . a a e e . f . he e . a e $|SSI| = 1000$, $r_n = 50\%$, $|SQFI| = 10$, $|t| = 12$, $N = |D| = 40000$,

Table 4. Comparison of accuracies

dataset	num. of records	num. of attributes(numeric)	num. of classes	C4.5	CBA	CMAR	CAEP	LSC-CAEP [comp. time (sec)]
Cleve	303	13(5)	2	.782	.828	.822	.833	.789 [38]
Ecoli	336	8(7)	8	.824	-	-	-	.831 [22]
Heart	270	13(6)	2	.808	.819	.822	.837	.845 [87]
Hepatitis	155	19(6)	2	.806	.818	.805	.830	.852 [26]
Iris	150	4(4)	3	.953	.947	.940	.947	.967 [0.1]
Glass	214	9(9)	7	.687	.739	.701	-	.681 [19]
Labor	57	16(8)	2	.793	.863	.897	-	.943 [0.1]
Wine	178	13(13)	3	.927	.950	.950	.971	.972 [52]
Zoo	101	16(0)	7	.922	.968	.971	-	.911 [19]

support=19: {class:good, duration-years:[2,2], working-hours:[33,40], wage-inc.-2nd-year(%):[4.0,5.8]}.

support=16: {class:good, duration-years:[3,3], working-hours:[35,40], wage-inc.-2nd-year(%):[3.5,5.0]}.

The e QFI gge a a . . . 1. ha he 1 c.ea e f . b. ab11 f . . 2 ea . . 3 ea . ba a ce 1 h ad 1 1 g . igh . . ge . . 1 g h . . a d 0.5% ~ 0.8% e . age 1 c.ea e. The f . 1 g . CAR ha 1 g high g . h , a e a e f . d 1 I 1 da a. The 1 igh . . he . ecie f 1 1 ca be ea . ed. growth rate=4.5: petal width:[1.4-2.5] → class:virginica growth rate=1.9: sepal length:[4.9-7.0], sepal width:[2.0-3.4] → class:setosa. The . . e g a . a 1 f he 1 e . a b . da 1 e he . . he 1 e . e a 1 . .

5 Discussion and Conclusion

The chec he a icab11 f LSC-CAEP . age da a e . , LSC-CAEP, C4.5 a d CBA e e a ed . Ce . . -I c . e da a c . a 1 g 199523 1 . a ce a d 40 a . lb e (. . e ic : 7 ca eg . ica : 33) 1 UCI KDD A . ch i e , a d c . . . ed ha he acc . ac a chie ed b LSC-CAEP 1 92.4% , hich 1 c . . a ab e 1 h 94.3% a d 94.0% f C4.5 a d CBA e . ec i e . F , he 1 . . e e . f he e f . a ce f LSC-CAEP 1 be add e . ed 1 f . e . d .

The . . e e . 1 e a . 1 LSC-CAEP a e he . . hich 1 $O(N \cdot g \cdot N)$ a d QFI-C . . . f Fig.1. The . a 1 a de 1 -c . . ec ed e . . e e . . e ic a . lb e p a e ea 1 de 1 ed 1 . . e ca f TIDLs . temp 1 MDCS b . 1 g he . . 1 de 1 b 1 a he . . e 1 Fig.2. He ce 1 1 $O(N)$ a . a 1 . . The 1 e a 1 . . f he . e . . f . . (5) . (10) 1 QFI-C . . a 1 e e . 1 e . I he . . ca e , a 1 . a ce 1 e . ed 1 each . . a h f . . he edge f a , eg 1 . he e 1 . a ce a e , a ged 1 a e 1 dic . a e , a d he . . bec . e $O(N^2)$. He e e , 1 he . . 1 e ca e hich 1 a e . . e ia de 1 di . lb 1 . , a . . 1 . $0 < r < 1$ f he 1 . a ce 1 he a e age a e . e a 1 ed 1 each . . a h. The . . . 1 he b . he 1 e $r^m N$ bec . e e . ha minsup he e m 1 he . . bec . f . . a h. Th . minsup $\leq r^m N$, a d m 1 a . . d $O(\cdot g \cdot N)$. Acc . di g , he e . 1 e a g 1 h 1 e ec ed . be $O(N \cdot g \cdot N)$.

O a e ab ed e cie . b ace c . e 1 g . . he 1 . e f . . e ic a d ca eg . ica da a 1 a e e 1 e a g 1 h . F , he . e a . . a che f c . . e 1 g a d c a 1 ca 1 f . a ge da a e . ca be de e . ed a . g hi 1 e .

Acknowledgement. The authors wish to thank D. A. De Amorim, ISIR, Osaka University, for his helpful comments. This research was partially supported by the Japanese Society for the Promotion of Science (JSPS) Grant-in-Aid for Scientific Research (B), 16300045, 2005.

References

1. Liu, B., Hsu, W., Ma, Y.: Integrating classification and association rule mining. Proc. of Fourth International Conference on Knowledge Discovery and Data Mining (1998)
2. Li, W., Han, J., Pei, J.: Cmar: Accurate and efficient classification based on multiple class-association rules. Proc. of First IEEE International Conference on Data Mining (2001) 369–376
3. Dong, G., Zhang, X., Wong, L., Li, J.: Caep: Classification by aggregating emerging patterns. Proc. of Second International Conference on Discovery Science, Lecture Notes in Computer Science **1721** (1999) 30–42
4. Agrawal, R., Gehrke, J., Gunopulos, D., Raghavan, P.: Automatic subspace clustering of high dimensional data for data mining applications. Proc. of the 1998 ACM SIGMOD international conference on Management of data (1998) 94–105
5. Procopiuc, C.M., Jones, M., Agarwal, P.K., Murali, T.M.: A monte carlo algorithm for fast projective clustering. Proceedings of the 2002 ACM SIGMOD international conference on Management of data (2002) 418–427
6. Kailing, K., Kriegel, H.P., Kroger, P.: Density-connected subspace clustering for high-dimensional data. Proc. Fourth SIAM International Conference on Data Mining (SDM'04) (2004) 246–257
7. Srikant, R., Agrawal, R.: Mining quantitative association rules in large relational tables. Proc. of 1996 ACM SIGMOD Int. Conf. on Management of Data (1996) 1–12
8. Wang, K., Hock, S., Tay, W., Liu, B.: Interestingness-based interval merger for numeric association rules. Proc. of 4th Int. Conf. on Knowledge Discovery and Data Mining (KDD) (1998) 121–128
9. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules. Proc. of 20th Int. Conf. on Very Large Data Bases (VLDB) (1994) 487–499

An Incremental Algorithm for Mining Generators Representation

Lijun Xu and Kanglin Xie

Department of Computer Science and Engineering, Shanghai JiaoTong University, 282#,
No.1954 HuaShan Road, Shanghai, China, 200030
lijunxu@sjtu.edu.cn, xie-kl@cs.sjtu.edu.cn

Abstract. This paper presents an efficient algorithm for maintaining the generator representation in dynamic datasets. The generators representation is a kind of lossless, concise representation of the set of frequent itemsets. Furthermore, the algorithm utilizes a novel optimization based on generators borders for the first time in the literature. Generators borders are the borderline between frequent generators and other itemsets. New frequent generators can be generated through monitoring them. Experiments show that our algorithm is more efficient than previous solutions.

1 Introduction

Frequent itemsets mining [1] is an important subject in many data mining applications, such as the discovery of association rules, correlations, sequential rules and episodes. A lot of algorithms have been proposed for this domain. But most algorithms assume that all transactions are available prior to the execution of the algorithm. However, in most cases this assumption does not hold. Many datasets are updated with blocks of data at regular time intervals. Recognizing the importance of the problem, many researchers [2-6, 11, 12] have proposed their solutions and efficient algorithms. The first incremental frequent itemsets mining algorithm, FUP, was proposed by Cheung et al. [3]. FUP2 [4] algorithm, adapted from FUP, can simultaneously handle deletions and additions. Two algorithms both adopt a level-wise search strategy like Apriori algorithm [1] and use the previous result for guiding the update. Feldman et al. [6] and Thomas et al. [11] proposed two similar algorithms respectively. The main idea of two algorithms is to keep track of frequent itemsets and the negative border that contains the itemsets form the borderline between frequent itemsets and infrequent itemsets. New frequent itemsets can be found by monitoring the negative border. Ayan et al. [2] presented UWEP algorithm, which follows the approach of FUP2. UWEP prunes the itemsets that will become infrequent by a look-ahead pruning strategy. ZIGZAG algorithm [12] is enlightened by GenMax [7] algorithm, an algorithm for discovering maximal frequent itemsets. It incrementally computes maximal frequent itemsets combining previous knowledge. But it may scan the dataset again to compute support values of some frequent itemsets that are not maximal frequent itemsets. Chi et al. [5] proposed Moment algorithm, which uses an in-memory data structure to monitor frequent closed itemsets and the itemsets that form the boundary between the frequent closed itemsets and the rest of the itemsets. Moment handles

new transactions or deleted transactions one by one, which may cause frequent changes of the boundary and affect the performance of the algorithm.

In this paper we present an efficient algorithm, called GBorder2, to maintain the generators representation in dynamic datasets. The generators representation is a kind of lossless, concise representation of the set of frequent itemsets. The usage of the generators representation can significantly reduce the times of data scans and the number of candidates in that the generators representation can be orders of magnitude smaller than the set of all frequent itemsets. Moreover, to the best of our knowledge, the algorithm introduces a novel optimization utilizing generators borders for the first time. Generators borders are the borderline between frequent generators and other itemsets. This optimization provides significantly computational or I/O savings as new frequent generators can be generated through monitoring generators borders.

2 Problem Definition

Let I be a set of items. A subset $X \subseteq I$ is called an itemset. An itemset with k items is called k -itemset. Let D be a transactional database, where each transaction is a subset of I . The number of transactions in D is denoted by $|D|$. During each update, obsolete transactions are removed and new transactions are added. Let d^+ be the set of newly added transaction, d^- be the set of deleted transactions and N be the updated dataset, i.e. $N = (D - d^-) \cup d^+$.

The support value of an itemset X , $\text{Sup}(X)$, is the number of the transactions that contain X . An itemset is frequent if it satisfies the minimum support threshold (θ). Let F be the set of frequent itemsets, i.e. $F = \{X | \text{Sup}(X) \geq \theta |D|\}$.

An itemset is a generator if none of its proper subsets has the same support as it has. We denote the set of generators by G and the set of frequent generators by FG , i.e. $FG = F \cap G$. Negative generators border, GB^- , is defined as the set of infrequent generators whose proper subsets are frequent generators. Positive generators border, GB^+ , is defined as the set of frequent non-generators whose proper subsets are generators. The generators representation consists of two components: (a) FG enriched by the support value for each itemset $X \in FG$; (b) GB^- . The following lists two important conclusions. Please refer to [8, 9] for more details.

Theorem 1. $X \in G \rightarrow \forall S \subset X, S \in G; X \notin G \rightarrow \forall S \supset X, S \notin G$.

Theorem 2. Let $X \subseteq I$. If $\exists Z \in GB^-$ and $Z \subseteq X$, then $X \notin F$. Otherwise, $X \in F$ and $\text{Sup}(X) = \min(\{\text{Sup}(S) | S \in FG \wedge S \subseteq X\})$.

3 GBorder2 Algorithm

GBorder2 algorithm is enlightened by the idea of the negative border [6, 10, 11]. GBorder2 maintains two kinds of generators borders: GB^- and GB^+ . GB^- defines the borderline between frequent generators and infrequent generators, and GB^+ defines the borderline between frequent generators and frequent non-generators. Most itemsets do not change their status (from frequent to infrequent, from infrequent to frequent, from generator to non-generator or from non-generator to generator) when a

small number of new transactions are added or a small portion of the dataset is removed. If the itemset does not change its status, nothing needs to be done except for updating its support value. Otherwise, as we shall present, the changes must come through generators borders.

Theorem 3. Let ChangedGB be a set of itemsets that belong to FG in N and belong to GB^+ or GB^- in D . If X is a frequent generator in N and is not a frequent non-generator in D , then there exists a subset $Y \subseteq X$, $Y \in \text{ChangedGB}$.

Proof: There are two possible cases for X :

1. X is a frequent non-generator in D . Let Y be the smallest subset of X that is a frequent generator in N but a frequent non-generator in D . As Y has minimal size, all its proper subsets are frequent generators in D . Thus Y belongs to GB^+ in D and Y belongs to ChangedGB.
2. X is an infrequent itemset in D . Let Y be the smallest subset of X that is a frequent generator in N but an infrequent generator in D . As Y has minimal size, all its proper subsets are frequent generators in D . Thus Y belongs to GB^- in D and Y belongs to ChangedGB.

3.1 Algorithm Description

The pseudo-code for GBorder2 algorithm is given in Fig. 1. We assume that each itemset X that belongs to frequent generators or generators borders (OldFG, Old GB^- or Old GB^+) and its support value in D , $\text{sup}(X, D)$, are already known.

The approach starts by scanning d^+ , d^- and computing the support values of all itemsets of OldFG, Old GB^- and Old GB^+ in d^+ and d^- respectively (Lines 1-3). Since the addition of new transactions and the deletion of obsolete transactions, some itemsets of OldFG, Old GB^- or Old GB^+ may change their status. Thus the frequent generators and the generators borders are determined again (Lines 4-6). ChangedGB contains the new frequent generators that originally belong to the generators borders in D (Line 7). It is used to generate candidates in the later steps.

Next, the candidates are generated and tested level by level like the classical Apriori algorithm [1] (Lines 8-26). $(i+1)$ -candidates (C_{i+1}), is generated based on i -itemsets of ChangedGB (Changed GB_i), new i -generators calculated in the last while-loop steps (G_i), i -generators (New FG_i) (Line 12). For each candidate X , the algorithm first determines $\text{Sup}(X, d^+)$ and $\text{Sup}(X, d^-)$ by scanning d^+ and d^- (Line 14). Then there are two possible cases when $\text{Sup}(X, D)$ is calculated. If X is infrequent in D , the algorithm has to scan D and determines its support value (Line 15-16). Otherwise, its support value can be directly retrieved from OldFG according to Theorem 2 (Lines 17-18). Finally the qualified candidates are added into NewFG (Line 23), New GB^- (Line 21) or New GB^+ (Line 25) respectively after updating their support values.

The while-loop steps (Lines 10-26) are performed only if ChangedGB is not empty. Thus unnecessary computing and I/O requirements are avoided if there is no new generator generated. Furthermore, the number of candidates can be considerably reduced even though these steps are performed.

Input: OldFG, OldGB⁻, OldGB⁺, N (N=(D-d⁻)∪d⁺) and θ
Output: NewFG, NewGB⁻ and NewGB⁺

- 1) for $X \in \text{OldFG} \cup \text{OldGB}^- \cup \text{OldGB}^+$
- 2) Scan d^+ , d^- and calculate $\text{Sup}(X, d^+)$, $\text{Sup}(X, d^-)$
- 3) $\text{Sup}(X, N) = \text{Sup}(X, D) + \text{Sup}(X, d^+) - \text{Sup}(X, d^-)$
- 4) $\text{NewFG} = \{X \mid X \in \text{OldFG} \cup \text{OldGB}^- \cup \text{OldGB}^+ \wedge \text{Sup}(X, N) \geq \theta \mid N \mid \wedge \forall S \subset X, \text{Sup}(X, N) < \text{Sup}(S, N)\}$
- 5) $\text{NewGB}^- = \{X \mid X \in \text{OldFG} \cup \text{OldGB}^- \cup \text{OldGB}^+ \wedge \text{Sup}(X, N) < \theta \mid N \mid \wedge \forall S \subset X, S \in \text{NewFG} \wedge \forall S \subset X, \text{Sup}(X, N) < \text{Sup}(S, N)\}$
- 6) $\text{NewGB}^+ = \{X \mid X \in \text{OldFG} \cup \text{OldGB}^- \cup \text{OldGB}^+ \wedge \text{Sup}(X, N) \geq \theta \mid N \mid \wedge \forall S \subset X, S \in \text{NewFG} \wedge \exists S \subset X, \text{Sup}(X, N) = \text{Sup}(S, N)\}$
- 7) $\text{ChangedGB} = \{X \mid X \in \text{OldGB}^- \cup \text{OldGB}^+ \wedge X \in \text{NewFG}\}$
- 8) $n = \max(\{i \mid \text{ChangedGB}_i \neq \emptyset\})$,
- 9) $G_0 = \emptyset, i = 0$
- 10) while ($G_i \neq \emptyset \vee i \leq n$)
- 11) $G_{i+1} = \emptyset$
- 12) $C_{i+1} = \{X \mid |X| = i+1 \wedge \exists i\text{-subset } S \subset X, S \in \text{ChangedGB}_i \cup G_i \wedge \forall i\text{-subset } S \subset X, S \in \text{NewFG}_i \cup \text{ChangedGB}_i\}$
- 13) for $X \in C_{i+1}$
- 14) Scan d^+ , d^- and calculate $\text{Sup}(X, d^+)$, $\text{Sup}(X, d^-)$
- 15) if $\exists S \subset X \wedge S \in \text{OldGB}^-$ then
- 16) Scan D and calculate $\text{Sup}(X, D)$
- 17) else
- 18) $\text{Sup}(X, D) = \min\{\text{Sup}(S, D) \mid S \subset X \wedge S \in \text{OldFG}\}$
- 19) $\text{Sup}(X, N) = \text{Sup}(X, D) + \text{Sup}(X, d^+) - \text{Sup}(X, d^-)$
- 20) if $\text{Sup}(X, N) < \theta \mid N \mid$ then
- 21) Add X into NewGB^-
- 22) else if $\forall S \subset X, \text{Sup}(X, N) < \text{Sup}(S, N)$ then
- 23) Add X into G_{i+1}
- 24) else
- 25) Add X into NewGB^+
- 26) $\text{NewFG} = \text{NewFG} \cup G_{i+1}, i = i+1$

Fig. 1. GBorder2 Algorithm

3.2 Discussions

GBorder2 handles the general case for transaction insertions as well as deletions. For the add-only case ($d^+ \neq \emptyset$ and $d^- = \emptyset$) or the delete-only case ($d^+ = \emptyset$ and $d^- \neq \emptyset$), there exists some improvements on the implementation of the algorithm.

For the add-only case, as we shall present in Theorem 4, a generator in D is still a generator in N. Then we can optimize GBorder2 by modifying Line 4-6 in Fig.1. The changes are shown in Fig. 2.

Theorem 4. Let X be a generator in D. If $d^- = \emptyset$, i.e. $N = D \cup d^+$, then X is still a generator in N.

- 4) $NewFG = \{X \mid X \in OldFG \cup OldGB^- \wedge Sup(X, N) \geq \theta \mid N \mid\} \cup \{X \mid X \in OldGB^+ \wedge Sup(X, N) \geq \theta \mid N \mid \wedge \forall S \subset X, Sup(X, N) < Sup(S, N)\}$
- 5) $NewGB^- = \{X \mid X \in OldFG \cup OldGB^- \wedge Sup(X, N) < \theta \mid N \mid \wedge \forall S \subset X, S \in NewFG\} \cup \{X \mid X \in OldGB^+ \wedge Sup(X, N) < \theta \mid N \mid \wedge \forall S \subset X, S \in NewFG \wedge \forall S \subset X, Sup(X, N) < Sup(S, N)\}$
- 6) $NewGB^+ = \{X \mid X \in OldGB^+ \wedge Sup(X, N) \geq \theta \mid N \mid \wedge \forall S \subset X, S \in NewFG \wedge \exists S \subset X, Sup(X, N) = Sup(S, N)\}$

Fig. 2. Optimizations for add-only case

Proof. Let S be an arbitrary subset of X. According the definition of generators, $Sup(X, D) < Sup(S, D)$. AS S is a subset of X, $Sup(X, d^+) \leq Sup(S, d^+)$. Then $Sup(X, N) = Sup(X, D) + Sup(X, d^+) < Sup(S, D) + Sup(S, d^+) = Sup(S, N)$.

So X is a generator in N.

For the delete-only case, a non-generator in D is still a non-generator in N (See Theorem 5). So any new generator must be infrequent generator in D. We have two improvements over the pseudo-code of GBorder2. The first one is presented in Fig. 3. The second one is that Lines 15-18 are replaced with Line 16 as none of the candidates are frequent in D.

- 4) $NewFG = \{X \mid X \in OldFG \cup OldGB^- \wedge Sup(X, N) \geq \theta \mid N \mid \wedge \forall S \subset X, Sup(X, N) < Sup(S, N)\}$
- 5) $NewGB^- = \{X \mid X \in OldFG \cup OldGB^- \wedge Sup(X, N) < \theta \mid N \mid \wedge \forall S \subset X, S \in NewFG \wedge \forall S \subset X, Sup(X, N) < Sup(S, N)\}$
- 6) $NewGB^+ = \{X \mid X \in OldFG \cup OldGB^- \wedge Sup(X, N) \geq \theta \mid N \mid \wedge \forall S \subset X, S \in NewFG \wedge \exists S \subset X, Sup(X, N) = Sup(S, N)\} \cup \{X \mid X \in OldGB^+ \wedge Sup(X, N) \geq \theta \mid N \mid \wedge \forall S \subset X, S \in NewFG\}$
- 7) $ChangedGB = \{X \mid X \in OldGB^- \wedge X \in NewFG\}$

Fig. 3. Optimizations for delete-only case

Theorem 5. Let X be a non-generator in D. If $d^+ = \emptyset$, i.e. $N = D - d^-$, then X is still a non-generator in N.

Proof. Let S be an proper subset of X and $Sup(X, D) = Sup(S, D)$. Obviously, any transaction in D that contains S also contain X. d^- is a portion of D and thus $Sup(X, d^-) = Sup(S, d^-)$. Then $Sup(X, N) = Sup(X, D) - Sup(X, d^-) = Sup(S, D) - Sup(S, d^-) = Sup(S, N)$. So X is a non-generator in N.

4 Experimental Results

We performed extensive experiments to evaluate GBorder2 algorithm. We compared it with FUP2 algorithm. We implemented two algorithms using Microsoft Visual C++ 6.0. We used the same data structures and subroutines in order to minimize any performance differences caused by minor differences in implementation. The two

algorithms are not fully optimized due to the time limitation. They were performed on a Pentium 1.2G processor with 1G MB, running Windows 2000.

We choose four datasets for the performance tests, which are publicly available from IBM Almaden Research Center (www.almaden.ibm.com/cs/quest/demos.html). The T10I4D100K dataset and the T40I10D100K dataset are synthetic datasets, while the connect dataset and the gazelle dataset are real-world datasets. Their characteristics are shown in Table 1.

Table 1. Characteristics of four datasets

Dataset	#Items	#Trans.	Avg. Trans. Len.	Max. Trans. Len.
T10I4D100K	1,000	100,000	3.7	31
T40I10D100K	1,000	100,000	8.5	77
gazelle	498	59,601	2.5	267
connect	130	67,557	43	43

We first conducted several experiments to evaluate the speed up of GBorder2 over FUP2. Without loss of generality, let $|D|=100K$ and $|d^+|=|d^-|=10K$. We duplicated and randomized each original dataset to obtain 110K transactions. Fig. 4 shows the results over different datasets. There are two interesting trends we observe:

1) For synthetic datasets, GBorder2 shows better performance for high support thresholds than low support thresholds. The reason is that the probability of generators borders expanding is higher at low support thresholds and as a result GBorder2 may have to scan the whole dataset.

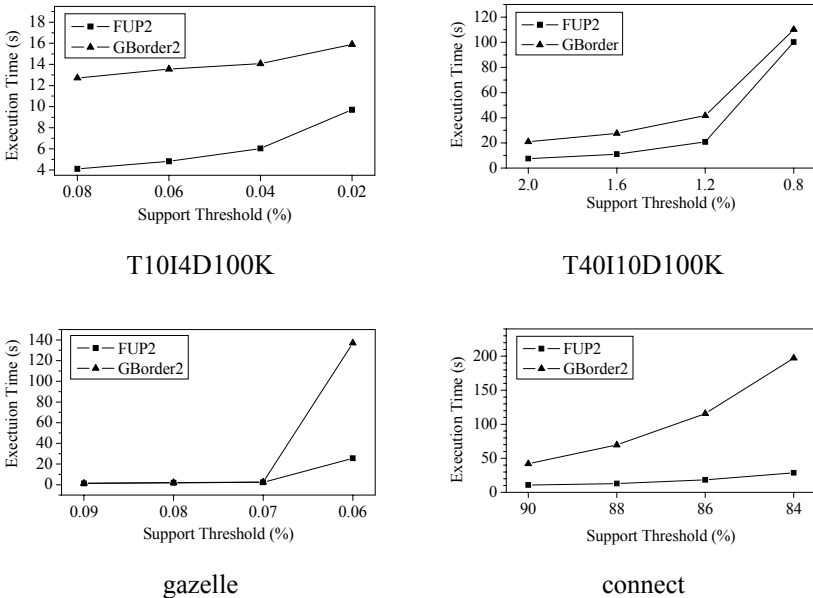


Fig. 4. Performance experiments

2) For real-world datasets, GBorder2 outperforms FUP2 throughout the entire range. Moreover, the performance gain of GBorder2 is larger for higher support thresholds. The phenomenon should be caused by the characteristics of real-world datasets. Real-world datasets are always strongly correlated datasets and a large number of frequent itemsets are non-generators for them. On the contrary, most frequent itemsets are generator for synthetic datasets.

Next, we conducted some experiments to find out if GBorder2 is able to deal with large datasets. Let $|D|=x$ and $|d^+|=|d^-|=x/10$, where x is varied in the experiments. We used a support threshold of 0.02% for the T10I4D100K dataset and 0.06% for the gazelle dataset. The results are plotted in Fig. 5. Obviously, the execution time of GBorder2 increases linearly as x increase, which implies that GBorder2 can handle large datasets well.

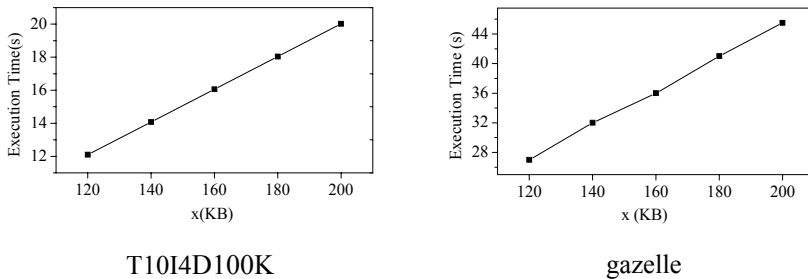


Fig. 5. Scale-up experiments

5 Conclusion

The paper focuses on the problem of frequent itemsets mining in dynamic datasets. Unlike existing incremental approaches, we propose an efficient algorithm to discover the generators representation using generators borders. The generators representation is a lossless, concise representation of frequent itemsets. New frequent generators can be computed by monitoring generators borders alone. To the best of our knowledge, it is the first incremental approach that combines the border technique and the generators representation. The usage of two techniques provides significant computational and I/O savings. Extensive experimental results show the efficiency of our approach.

A number of lossless concise representations have been proposed [8]. All these representations, except for frequent closed itemsets, consist of two components: one main component and several borders. All border representations, except for the generators representation, are about two orders of magnitude more concise than frequent closed itemsets in practice. Due to the common characteristics of all border representations, our algorithm can be extended to update other border representations in an incremental manner.

Reference

1. R. Agrawal and R. Srikant. Fast algorithms for mining association rules. In Proc. of VLDB'94 (1994) 487-499.
2. N. Ayan, A. Tansel and E. Arkun. An efficient algorithm to update large itemsets with early pruning. In Proc. of the 5th ACM-SIGKDD (1999) 287-291.
3. D. Cheung, J. Han, V. Ng and C. Y. Wong. Maintenance of discovered association rules in large databases: An incremental updating technique. In Proc. of the 12th Data Engineering (1996) 106-114.
4. D. Cheung, S. Lee, and B. Kao. A general incremental technique for maintaining discovered association rules. In Proc. of the 5th Database Systems for Advanced Applications (1997) 1-4.
5. Y. Chi, H. Wang, P. S. Yu and R. R. Muntz. Moment: maintaining closed frequent itemsets over a stream sliding window. In Proc. of ICDM'04 (2004) 59-66.
6. R. Feldman, Y. Aumann, A. Amir, and H. Mannila. Efficient algorithms for discovering frequent sets in incremental databases. In Proc. of ACM SIGMOD Workshop on Research Issues on Data Mining and Knowledge Discovery (1997) 59-66.
7. K. Gouda and M. Zaki. Efficiently mining maximal frequent itemsets. In Proc. of ICDM'01 (2001) 163-170.
8. M. Kryszkiewicz. Concise representations of association rules. In Proc. of the ESF Exploratory Workshop on Pattern Detection and Discovery (2002) 92-109.
9. N. Pasquier, Y. Bastide, R. Taouil and L. Lakhal. Pruning closed itemset lattices for association rules. In Proc. of BDA'98 (1998) 177-196.
10. H. Toivonen. Sampling large databases for association rules. In Proc. of VLDB'96 (1996) 134-145.
11. S. Thomas, S. Bodagala, K. Alsabti, and S. Ranka. An efficient algorithm for the incremental updation of association rules. In Proc. of KDD'97 (1997) 263-266.
12. A. Veloso, W. Meira Jr., M. B. de Carvalho, B. Póssas, S. Parthasarathy, and M. Zaki. Mining frequent itemsets in evolving databases. In Proc. of SIAM'02 (2002).

Hybrid Technique for Artificial Neural Network Architecture and Weight Optimization

Ceber Zaçetti and Teresa Berarda Ludermir

Center for Informatics, Federal University of Pernambuco (UFPE),
P.O. Box 7851, 50.732-970, Recife, PE, Brazil
{cz, tbl}@cin.ufpe.br

Abstract. This work presents a technique that integrates the heuristics tabu search, simulated annealing, genetic algorithms and backpropagation. This approach obtained promising results in the simultaneous optimization of the artificial neural network architecture and weights.

1 Introduction

Optimization is the process of finding the best solution for a problem from a group of possible solutions. An optimization problem has a objective function and a group of restrictions, both related to the decision variables of the problem. Genetic Algorithms (AG) [3], Simulated Annealing (SA) [1] and Tabu Search (TS) [2] are iterative algorithms used to solve difficult combinatorial optimization problems. These three algorithms are the most popular from a class of optimization algorithms known as generative algorithms. A three optimization heuristics have similarities [4]: (1) They are approximate (heuristic) algorithms, i.e., they do not assure the finding of a local optimum solution; (2) They are blind that they do not know when they have reached a optimum solution, and therefore, must be terminated step; (3) They have a high combinatorial growth, i.e., they occasionally accept uphill (bad) moves; (4) They are generative, i.e., they can easily be generated to implement a combinatorial optimization problem; and that is required is that they have a suitable solution representation, a cost function, and a mechanism to traverse the search space; and (5) Under certain conditions, they asymptotically converge to a optimum solution.

This paper presents a new technique that integrates the main particularities of these three heuristics. This technique is evaluated in the simultaneous optimization of the number of connections and eight connection weights among processing units of the Multi-Layer Perceptron neural network (MLP) [5].

The MLP trained by the backpropagation algorithm (BP) is one of the most used connectionist models in the literature. To obtain successful use, the network topology parameter is important. A lack of connections can reduce the network capabilities of solving the investigated problem as a result of the inadequacy of adjustable parameters, whereas an excess of connections can cause overfitting to the training data and fail to have an adequate generalization capacity. In general, the training of the MLP neural networks is accomplished through successive attempts

with direct, iterative, population genetic algorithms usually reaching satisfactory results for the problem. Besides computing time, this process can establish effective architectures with unnecessary connections and nodes. Moreover, the larger the population, the more complete the available adjustment of these connections becomes. Thus, the simultaneous use of optimization of architectures and weights of artificial neural networks is a more interesting approach to the generation of efficient networks with small populations.

2 Search Heuristics Description

The genetic algorithm is characterized by a parallel search of the state space as against a typical depth search through complete optimization techniques. The parallel search is accomplished by keeping a set of populations useful for the optimization problem, called populations. An individual in the population is a string of symbols and is an abstract representation of the solution. The symbols are called genes and each string of genes is termed a chromosome. The individuals in the population are evaluated through a fitness measure. The populations of chromosomes evolve from the generation to the next through the use of different genetic operators: (1) unary operators, such as mutation and inversion, which alter the genetic structure of a single chromosome; and (2) higher-order operators, referred to as crossover, which consists of binary genetic material from two selected parent chromosomes. The parent chromosomes are chosen by a selection technique [3].

In the experiments performed, each chromosome is represented as described in Section 3. The initial population is defined with a size of 10 chromosomes. The chromosomes are classified by Rank Based Fitness Scaling [8]. The parents chosen for the next generation is accomplished by a probabilistic manner, using Uniform Stochastic Sampling [8]. Elitism is also used, with a probability of 10%. For the combination of the parent chromosomes, the crossover operator Uniform Crossover [9] is used, with a probability of 80%. The mutation operator used is the Gaussian Mutation [6], with a probability of 10%. The stopping criteria are: (1) the GL_5 criterion, this criterion provides an idea of the generalization of the training data and it is sufficient useful and overfitting. It is defined as the increase in the validation error rate to the minimum validation error; and (2) a maximum number of 500 generations.

The simulated annealing method is different of the others search methods in that uphill moves are occasionally accepted to escape from a minimum. The search process consists of a sequence of iterations. Each iteration consists of random change of the current solution to create a new solution in its neighborhood. Once a new solution is created, the corresponding change in the cost function is computed to decide if the new solution can be accepted. If the new solution cost is lower than the current solution cost, is accepted. Otherwise, the Metropolis's criterion is verified [10], based on the Boltzmann probability. A random number $d \in [0, 1]$ is generated from a uniform distribution. If $d \leq e^{-\frac{\Delta E}{T}}$, here ΔE is the change in the cost function and T is a parameter called temperature, then the new solution is accepted as the current solution. If not, the current solution is unchanged and the process continues from the current solution.

The algorithm as rigidly derived from thermodynamic simulations. Thus, the parameter T is referred as temperature and the temperature reduction process is called the cooling process. The chosen cooling strategy as follows. According to this rule, the new temperature is equal to the current temperature multiplied by a temperature factor (smaller than one, but close to one) [11]. The initial temperature is set to 1, and the temperature factor is set to 0.9. The temperature is decreased at each 10 iterations, with a maximum number of 1,000 iterations. The stop criterion GL_5 as used.

Tabu search is an iterative search algorithm characterized by the use of a forbidden memory. In this method, each iteration consists of the evaluation of a certain amount of solutions (neighboring moves). The best of these solutions (in terms of cost function) is accepted. However, the best candidate solution may not improve the current solution. Thus, the algorithm chooses the solution that produces the largest improvement or the smallest deterioration in the cost function. This strategy allows the method to escape from local minima. A tabu list is used to store a certain amount of recently visited solutions. The solutions in tabu list are marked as forbidden to subsequent iterations. The tabu list registers T visited solutions. When the list is full, a new move is registered and substituted to the oldest solution kept in the list.

In the present work, a neighborhood with 20 solutions is used, and the algorithm chooses the best T -tabu solution. The primary criterion [6] as used to compare solutions. A solution is considered ideologically the tabu solution if: (1) each correct bit in the solution is ideologically the corresponding correct bit in the tabu solution; and (2) each correct bit in the solution is within $\pm N$ of the corresponding correct bit in the tabu solution. The parameter N is a real number with a value of 0.001. A maximum number of 100 iterations is allowed. The stop criterion employed as a GL_5 .

3 Integration of Simulated Annealing, Tabu Search and Genetic Algorithms

The simulated annealing method has the ability to escape from local minima through the choice between accepting or discarding a new solution that creates cost (uphill moves). The tabu search method, in contrast, evaluates the group of solutions at each iteration (instead of only one solution as in simulated annealing). This makes a tabu search faster, as it generally needs less iterations to converge. The genetic algorithm evolves, in turn, in a sequence of iterations, where a group of solutions evolves through selection processes and reproduction. This process, which is more elaborate than the other algorithms, can result in solutions with a better quality.

These observations motivated the proposal of a optimization technique (GaTSa) that combines the main particularities of genetic algorithms, simulated annealing and tabu search in a alternative and their imitations. In general terms: at each iteration, a group of solutions is generated, starting from the micro-evolution of the current population, as in genetic algorithms. The cost of

each solution is evaluated, and the best solution is chosen, as in tabu search. However, different from a tabu search, this solution is not always accepted. The acceptance criterion is the same used in the simulated annealing algorithm - if the chosen solution has a smaller cost than the current solution, it is accepted; otherwise, it can either be accepted or not, depending on a probabilistic calculation. This probability is given by the same expression used in the simulated annealing method. The visited solutions are marked as tabu, as in a tabu search. During the optimization process, the best solution found is stored, that is, the final solution comes back through the method.

Algorithm 1. Pseudocode algorithm Pseudocode

1. $P_0 \leftarrow$ initial population with K solutions s_k
 2. $T_0 \leftarrow$ initial temperature
 3. $I_T \leftarrow$ iterations number
 4. Update S_{BSF} with s_k of the P_0 (best solution found so far)
 5. For $i = 0$ to $I_{max} - 1$
 6. If $i + 1$ is not a multiple of I_T
 7. $T_{i+1} \leftarrow T_i$
 8. Else
 9. $T_{i+1} \leftarrow$ new temperature
 10. If validation based stopping criteria are not satisfied
 11. Stop global search execution
 12. For $j = 0$ to g_n
 13. Generate a new population P' from P_i
 14. $P_i \leftarrow P'$
 15. Choose the best solution s_k from P_i
 16. If $f(s') < f(s_k)$
 17. $s_{k+1} \leftarrow s'$
 18. Else
 19. $s_{k+1} \leftarrow s'$ with probability $e^{\frac{f(s') - f(s_k)}{T_{i+1}}}$
 20. If $f(s_{k+1}) < f(S_{BSF})$
 21. Update S_{BSF}
 22. End For
 23. Keep the topology contained in S_{BSF} constant and use the weights as initial ones for training with the backpropagation algorithm
-

The pseudocode algorithm pseudocode is presented in Algorithm 1. Let S be a group of solutions and a real cost function, the pseudocode algorithm searches the global minimum s , such that $f(s) \leq f(s'), \forall s' \in S$. The process finishes after I_{max} iterations or if the stopping criterion based on the validation error is satisfied. The best found solution S_{BSF} (. . .) is returned. The cooling process updates the temperature T_i of the iteration i at each I_T algorithm iterations. At each iteration, a population with k solutions is generated. A genetic micro-evolution of g_n generations is used to generate this population from the current population. Moreover, at the end of the global search (GaTSa), a hybrid training is used, combining the pseudocode with a local search technique. The local search technique can be implemented, for instance, by the backpropagation algorithm.

Each solution is coded in a vector. This vector represents the connections among the processing units of the MLP artificial neural network. Each of these connections is specified by two parameters: (a) the connection bit, a binary value that symbolizes the existence or absence of a connection; and (b) the connection weight, which is a real number. If the connection bit is equal to zero,

its associated eight is considered, for the connectivity degrees, the first is the network. A possible connectivity sampling adapters are considered.

Directed from the connectivity algorithms that generate essential at the end of the process, iterative algorithms rigidity possible (candidate) solutions at each iteration. The cost function is used to evaluate the performance among consecutive iterations and select the solution that minimizes (or maximizes) a objective function.

The cost function for the investigated problem is the arithmetic average between: (1) the classification error of the training set (percentage of incorrect classified training patterns); and (2) the percentage of connectivity used by the artificial neural network. Therefore, the algorithms try to minimize both network performance and processing complexity. Obviously, networks (i.e., networks with at least one unit in the hidden layer) are considered.

The operator for the generation of neighbors is used to derive essential from the current solution. The method used in simulation is defined as follows: (1) the connectivity bits for the current solution are changed according to a given probability, which the present network is set to 20%. This operator defines a new network connectivity and creates new ones. Next, a random number taken from a uniform distribution in $[-1; +1]$ is added to each connectivity weight. These two steps can change both the positive and negative connectivity weights to produce a new neighbor solution.

4 Experiments and Results

Real data is used in the experiments. The problem aims to classify disorder patterns obtained through an artificial sense. The data characteristics analyzed are from three different stages (years 1995, 1996 and 1997) of the same commercial wine (A madm, Brazil) produced with merlot-type grapes. The artificial sense used is composed of six distinct conductive polymer sensors constructed with a electronic chemical deposition of polypyrrole using different deposition rates. Three data acquisitions were performed. In each acquisition for each individual stage, the resistance value of each sensor was recorded for five seconds. A set of six values from the six sensors at the same time was considered a pattern. Thus, each acquisition contains 1,800 patterns (600 from each stage). There were three acquisitions and 5,400 patterns of data.

In previous works with this database, the best performance obtained by the MLP was achieved by a architecture with 6 processing units in the input layer, 4 processing units in the hidden layer and 3 processing units in the output layer [7]. This topology was kept constant as the maximum architecture in the optimization experiments performed. In a investigated algorithms, the parameter configurations were maintained at the standard configurations readjusted based on previous experiments. The values used may not be the best values for the problem, but the objective of the present paper is to demonstrate the potentialities of the techniques and not the idea of algorithms configurations.

Table 1 presents the average performance of each investigated optimization technique. These results were obtained for each technique in the optimization

of the number of connections and eight connections of a MLP artificial neural network. The parameters evaluated were: (1) Squared Error Percentage (SEP) and the classification error (Class) of training, validation and test sets; (2) algorithm iteration number; (3) artificial neural network number; and (4) the temperature value. The following table displays the average results of 10 simulations. Each simulation consists 30 different runs of the algorithm.

Table 1. Optimization techniques performance

Technique	Training		Validation		Test		Iterations	Connections	Temperature
	SEP	Class	SEP	Class	SEP	Class			
TS	18,74	5,44	18,86	5,88	18,75	5,3805	51	11,42	-
SA	19,65	6,91	19,76	7,47	19,65	6,9331	715	11,77	0,0085
GA	21,66	15,88	21,73	16,52	21,66	15,9240	315	16,64	-
GaTSa	18,69	3,58	18,76	3,81	18,69	3,5664	46	8,33	0,7098
GaTSa + BP	4,78	-	2,41	-	2,14	2,8684	86	8,33	0,7098
BP	6,25	-	3,15	-	2,84	6,7854	90	36	-

The technique that combines the heuristics of tabu search, simulated annealing and genetic algorithms obtained the best result performance. This technique was better even without using the local search heuristic to optimize the artificial neural network connections. The average classification error obtained was 2.87%, with an average of 8 connections from 36 possible connections in a fully connected neural network. Using a fully connected network, the local optimization technique backpropagation obtained an average error of 6.78%.

The genetic algorithms, tabu search and simulated annealing methods incorporate domain specific knowledge in their local search heuristics. The stochastic elements of stochastic-determinism, which helps the search escape from local minima. Therefore, the use of a suitable stochasticity that provides feedback to the algorithm as the search progresses. The main difference among them is how and where domain-specific knowledge is used. Furthermore, simulated annealing such knowledge is mainly included in the stochasticity. Simulated annealing perturbations are selected randomly, and perturbations are accepted or rejected according to a probability.

In the case of genetic algorithms, domain specific knowledge is employed in all phases. The fitness function values, the reproductivity selection, genetic operators, as well as the generation of the population, incorporate domain-specific knowledge. Tabu search is different from the above heuristics in that it has a explicit memory component. At each iteration, the neighborhood of the current solution is partially explored, and a move is made to the best non-tabu solution in that neighborhood. The neighborhood function, together with the size and content of the tabu list, is problem specific. The direction of the search is assisted by memory structures.

The proposed integration uses a larger amount of information in the problem domain and uses this information in practical search phases. This is possible through the integration of the main particularities of the three investigated search

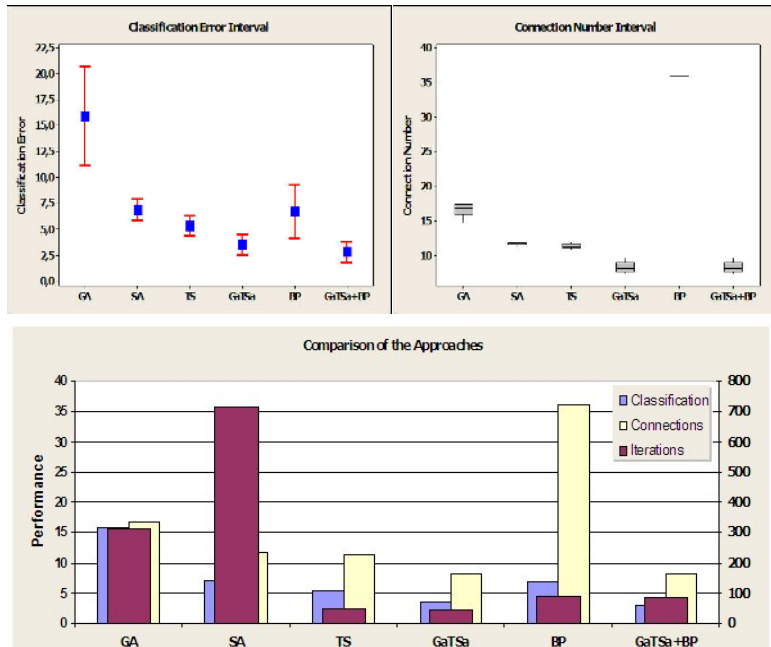


Fig. 1. Result analysis

heuristics. Moreover, the proposed technique has three well-defined stages: a global search phase, where it makes use of the capacity for generating solutions for the genetic algorithms, the cloning process and connectivity of the simulated annealing as well as the memory characteristics of the tabu search technique; a local search phase, where it makes use of characteristics such as gradient descent for a more precise solution adjustment. These characteristics can obtain better solutions for the investigated problems, with a short search time, low computational cost and minimize the investigated search space.

Figure 1 presents graphs comparing the performance of the investigated techniques. The proposed technique obtained the best results regarding the classification error, final network connectivity, number and the number of iterations needed for architecture optimization.

5 Final Remarks

This work presented a technique that integrates the heuristics of tabu search, simulated annealing, genetic algorithms and backpropagation. In the simulated annealing optimization of the connectivity, number and connectivity values of the Multilayer Perceptron neural network, this technique obtained promising results in comparison with the isolated techniques. The proposed technique combines strategies of global and local searches, presenting promising results regarding

the investigated search space, computational cost and search time. The investigated problem defines a critical subject, the stability versus plasticity relation in the training of artificial neural networks.

With out a deeper investigation, it is not possible to say if these results can be extended to other problem classes. An interesting theoretical study produced a number of theoretical results that the average performance of a pair of iterative (deterministic-recurrent-deterministic) algorithms across all problems is identical. Thus, if a algorithm performs well on a certain class of investigated problems, it is necessary to pass for that with degraded performance on the remaining set of problems [12]. Future investigations should consider this presupposition and verify the performance of this optimization technique on other problems.

Acknowledgments

The authors would like to thank CNPq, CAPES and FINEP (Brazilian research agencies) for their financial support.

References

1. Kirkpatrick, S., Gellat Jr, C. D., Vecchi, M. P.: Optimization by simulated annealing. *Science* **220** (1983) 671–680
2. Glover, F.: Future paths for integer programming and links to artificial intelligence. *Computers and Operation Research* **13** (1986) 533–549
3. Goldberg, D. E.: *Genetic Algorithms in Search, Optimization and Machine Learning*. Addison-Wesley Longman Publishing Co., Inc. (1989) 372 pages
4. Sait, S. M., Youssef, H.: *Iterative Computer Algorithms with Applications in Engineering: Solving Combinatorial Optimization Problems*. IEEE Computer Society Press (1999) 387 pages
5. Rumelhart, D. E., Hilton, G. E., Williams, R. J.: Learning Representations by Backpropagation Errors. *Nature* **323** (1986) 533–536
6. Sexton, R. J., Alidaee, B., Dorsey, R. E., Johnson, J. D.: Global Optimization for Artificial Neural Networks: A Tabu Search Application. *European Journal of Operational Research* **106:2-3** (1998) 570–584
7. Yamazaki, A., de Souto, M.C.P., Ludermir, T.B.: Optimization of Neural Network Weights and Architectures for Odor Recognition using Simulated Annealing. In *Proceedings International Joint Conference on Neural Networks* (2002) 547–552
8. Baker, J. E.: Reducing bias and inefficiency in the selection algorithm. In *Proceedings of the Second International Conference on Genetic Algorithms and their application*. Lawrence Erlbaum Associates (1987) 14–21
9. Sywerda, G.: Uniform crossover in genetic algorithms. In *Proceedings of international conf. on Genetic algorithms*. Morgan Kaufmann Publishers Inc. (1989) 2–9
10. Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., Teller, E.: Equation of state calculations by fast computing machines. *The Journal of Chemical Physics* **21:6** 1087–1092
11. Pham, D. T., Karaboga, D.: *Intelligent Optimisation Techniques*. Springer-Verlag New York (1998) 312 pages
12. Wolpert, D. H., Macready, W. G.: No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation* **1:1** (1997) 67–82

Author Index

- Abu-Mostafa, Yaser S. 157
Albertí, Pere 462
Appice, Annalisa 169
Atzori, Maurizio 10
Avesani, Paolo 343
Azevedo, Paulo J. 96
- Bahamonde, Antonio 462
Barajas, Jorge Mario 429
Beerenwinkel, Niko 285
Berthold, Michael R. 1
Bistarelli, Stefano 22
Bodon, Ferenc 437
Bonchi, Francesco 10, 22
Borges, José 34
Boulicaut, Jean-François 651
Bringmann, Björn 46
- Cai, Deng 445
Cardie, Claire 2
Carvalho, Deborah R. 453
Castelli, Vittorio 355
Ceccherini-Silberstein, Francesca 285
Ceci, Michelangelo 169
Chakrabarti, Deepayan 133
Chakraborti, Sutanu 380
Chen, Shyh-Kwei 368
Choong, Yeow Wei 205
Chowdhury, Abdur 561
Cunningham, Pádraig 486
- Däumer, Martin 285
Davidson, Ian 59
De Raedt, Luc 3
de Souza, Jerffeson Teixeira 667
del Coz, Juan José 462
Degenhard, Andreas 331
Demichelis, Francesca 343
Díez, Jorge 462
Ding, Chris 71
Domingos, Pedro 297
Dong, Lin 84
Du, Wenliang 643
- Ebecken, Nelson 453
Esfandiari, Babak 634
Ester, Martin 527
- Faloutsos, Christos 133
Fätkenheuer, Gert 285
Feldman, Ronen 217
Ferreira, Pedro Gabriel 96
Fischer, Ingrid 392
Frank, Eibe 84, 240, 675
Freitas, Alex A. 453
Fresko, Moshe 217
- Gábor, Bálint 470
Geurts, Pierre 478
Giannotti, Fosca 10
Govaert, Gérard 609
Greene, Derek 486
Greiner, Russell 121
Guo, Jun 264
Gyenes, Viktor 470
- Hall, Mark 675
Han, Jiawei 404, 445, 527
He, Xiaofei 445
He, Xiaofeng 71
Hilario, Melanie 536
Ho, Eric K.Y. 544
Ho, Tu Bao 321, 617
Hoffmann, Daniel 285
Holmes, Geoffrey 495
Huang, Jin 503, 511
Huang, Shao-bing 601
- Ilin, Alexander 519
- Japkowicz, Nathalie 667
Jin, Wen 527
- Kaiser, Rolf 285
Kalousis, Alexandros 536
Keogh, Eamonn 6, 577
Kirkby, Richard 495
Kleinberg, Jon 133
Knobbe, Arno J. 544

- Knuteson, Bruce 552
 Kohavi, Ron 7
 Kolcz, Aleksander 561
 Korn, Klaus 285
 Koychev, Ivan 380
 Kramer, Stefan 84
 Kriegel, Hans-Peter 417
- Laasonen, Kari 569
 Lakaemper, Rolf 577
 Latecki, Longin Jan 577
 Laurent, Anne 205
 Laurent, Dominique 205
 Law, Yan-Nei 108
 Lee, Chang-Hwan 585
 Lee, Chi-Hoon 121
 Lengauer, Thomas 285
 Leskovec, Jurij 133
 Levene, Mark 34
 Li, Haiquan 146
 Li, Jinyan 146
 Li, Ling 157
 Li, Qunxia 264
 Li, Wenyuan 593
 Li, Xue 429
 Lin, Hsuan-Tien 157
 Ling, Charles X. 274, 503, 511
 Liu, Gang 264
 Lórinicz, András 470
 Lothian, Rob 380
 Ludermir, Teresa Bernarda 709
 Lv, Tian-yang 601
- Malerba, Donato 169
 Matias, Yossi 8
 Matwin, Stan 667
 Mavroeidis, Dimitrios 181
 McGinty, Lorraine 228
 Megalooikonomou, Vasilis 577
 Meinel, Thorsten 392
 Motoda, Hiroshi 692
- Nadif, Mohamed 609
 Nakanishi, Koutarou 692
 Nattkemper, Tim 331
 Ng, Wee-Keong 593
 Nguyen, Canh Hao 617
 Nguyen, Son N. 625
 Nock, Richard 634
- Oette, Mark 285
 Olivetti, Emanuele 343
 Ong, Kok-Leong 593
 Orłowska, Maria E. 625
- Paşca, Marius 193
 Pedreschi, Dino 10
 Pei, Jian 684
 Pensa, Ruggero G. 651
 Perno, Carlo-Federico 285
 Pfahringer, Bernhard 495
 Philippsen, Michael 392
 Plantevit, Marc 205
 Polat, Huseyin 643
 Prados, Julien 536
 Pratap, Amrit 157
- Ratanamahatana, Chotirat Ann 577
 Ravi, S.S. 59
 Reilly, James 228
 Rexhepaj, Elton 536
 Robardet, Céline 651
 Rockstroh, Jürgen K. 285
 Rosenfeld, Benjamin 217
- Salamó, Maria 228
 Sañudo, Carlos 462
 Satou, Kenji 321
 Savary, Lionel 659
 Schmidberger, Gabi 240
 Schmidt, Mark 121
 Schmidt-Thieme, Lars 437
 Schneider, Karl-Michael 252
 Shao, Zheng 445
 Shen, Haifeng 264
 Sheng, Shengli 274
 Sing, Tobias 285
 Singla, Parag 297
 Smyth, Barry 228
 Soh, Donny 146
 Sumner, Marc 675
 Svicher, Valentina 285
- Talia, Domenico 309
 Teisseire, Maguelonne 205
 Theobald, Martin 181
 Tran, Tuan Nam 321
 Tresp, Volker 417
 Trunfio, Paolo 309
 Tsatsaronis, George 181

- Vagena, Zografoula 355
Valpola, Harri 519
Varini, Claudio 331
Vazirgiannis, Michalis 181
Veeramachaneni, Sriharsha 343
Verta, Oreste 309
Vilalta, Ricardo 552
Vlachos, Michail 355, 368
- Walter, Hauke 285
Wang, Haixun 684
Wang, Qiang 577
Wang, Zheng-xuan 601
Washio, Takashi 692
Wehenkel, Louis 478
Weikum, Gerhard 181
Wiratunga, Nirmalie 380
Wong, Limsoon 146
- Wörlein, Marc 392
Wu, Kun-Lung 368
- Xie, Kanglin 701
Xing, Yu-hui 601
Xu, Lijun 701
- Yan, Xifeng 445
Yin, Xiaoxin 404
Yu, Kai 417
Yu, Philip S. 355, 368
Yu, Shipeng 417
- Zanchettin, Cleber 709
Zaniolo, Carlo 108
Zeitouni, Karine 659
Zimmermann, Albrecht 46
Zuo, Wan-li 601