

Commenced Publication in 1973

Founding and Former Series Editors:

Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

Editorial Board

David Hutchison

Lancaster University, UK

Takeo Kanade

Carnegie Mellon University, Pittsburgh, PA, USA

Josef Kittler

University of Surrey, Guildford, UK

Jon M. Kleinberg

Cornell University, Ithaca, NY, USA

Alfred Kobsa

University of California, Irvine, CA, USA

Friedemann Mattern

ETH Zurich, Switzerland

John C. Mitchell

Stanford University, CA, USA

Moni Naor

Weizmann Institute of Science, Rehovot, Israel

Oscar Nierstrasz

University of Bern, Switzerland

C. Pandu Rangan

Indian Institute of Technology, Madras, India

Bernhard Steffen

University of Dortmund, Germany

Madhu Sudan

Massachusetts Institute of Technology, MA, USA

Demetri Terzopoulos

University of California, Los Angeles, CA, USA

Doug Tygar

University of California, Berkeley, CA, USA

Gerhard Weikum

Max-Planck Institute of Computer Science, Saarbruecken, Germany

Simeon J. Simoff Michael H. Böhlen
Arturas Mazeika (Eds.)

Visual Data Mining

Theory, Techniques and Tools
for Visual Analytics



Springer

Volume Editors

Simeon J. Simoff
University of Western Sydney
School of Computing and Mathematics
NSW 1797, Australia
E-mail: s.simoff@uws.edu.au

Michael H. Böhlen
Arturas Mazeika
Free University of Bozen-Bolzano
Faculty of Computer Science
Dominikanerplatz 3, 39100 Bozen-Bolzano, Italy
E-mail: {boehlen,arturas}@inf.unibz.it

Library of Congress Control Number: 2008931578

CR Subject Classification (1998): H.2.8, I.3, H.5

LNCS Sublibrary: SL 3 – Information Systems and Application, incl.
Internet/Web and HCI

ISSN 0302-9743
ISBN-10 3-540-71079-5 Springer Berlin Heidelberg New York
ISBN-13 978-3-540-71079-0 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
springer.com

© Springer-Verlag Berlin Heidelberg 2008
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India
Printed on acid-free paper SPIN: 12280612 06/3180 5 4 3 2 1 0

Foreword

Visual Data Mining—Opening the Black Box

Knowledge discovery holds the promise of insight into large, otherwise opaque datasets. The nature of what makes a rule interesting to a user has been discussed widely¹ but most agree that it is a subjective quality based on the practical usefulness of the information. Being subjective, the user needs to provide feedback to the system and, as is the case for all systems, the sooner the feedback is given the quicker it can influence the behavior of the system.

There have been some impressive research activities over the past few years but the question to be asked is why is visual data mining only now being investigated commercially? Certainly, there have been arguments for visual data mining for a number of years – Ankerst and others² argued in 2002 that current (autonomous and opaque) analysis techniques are inefficient, as they fail to directly embed the user in dataset exploration and that a better solution involves the user and algorithm being more tightly coupled. Grinstein stated that the “*current state of the art data mining tools are automated, but the perfect data mining tool is interactive and highly participatory,*” while Han has suggested that the “*data selection and viewing of mining results should be fully interactive, the mining process should be more interactive than the current state of the art and embedded applications should be fairly automated*”². A good survey on techniques until 2003 was published by de Oliveira and Levkowitz³.

However, the deployment of visual data mining (VDM) techniques in commercial products remains low. There are, perhaps, four reasons for this. First, VDM, as a strong sub-discipline of data mining only really started around 2001. Certainly there was important research before then but as an identifiable sub-community of data mining, the area coalesced around 2001. Second, while things move fast in IT, VDM represents a shift in thinking away from a discipline that itself has yet to settle down commercially. Third, to fully contribute to VDM a researcher/systems developer must be proficient in both data mining and visualization. Since both of these are still developing themselves, the pool from which to find competent VDM researchers and developers is small. Finally, if the embedding is to be done properly, the overarching architecture of the knowledge

¹ Geng, L., Hamilton, H.J.: Interestingness measures for data mining: A survey. *ACM Computing Surveys* 38 (2006)

² Ankerst, M.: The perfect data mining tool: Automated or interactive? In: Panel at ACM SIGKDD 2002, Edmonton, Canada. ACM, New York (2002)

³ de Oliveira, M.C.F., Levkowitz, H.: From visual data exploration to visual data mining: a survey. *IEEE Transactions on Visualization and Computer Graphics* 9, 378–394 (2003)

discovery process must be changed. The iterative paradigm of *mine and visualize* must be replaced with the data mining equivalent of direct manipulation⁴.

Embedding the user within the discovery process, by, for example, enabling the user to change the mining constraints, results in a finer-grained framework as the interaction between user and system now occurs *during* analysis instead of *between* analysis runs. This overcomes the computer's inability to incorporate evolving knowledge regarding the problem domain and user objectives, not only facilitating the production of a higher quality model, but also reducing analysis time for two reasons. First, the guidance reduces the search space at an earlier stage by discarding areas that are not of interest. Second, it reduces the number of iterations required. It also, through the Hawthorn Effect, has the effect of improving the user's confidence in, and ownership of, the results that are produced⁵.

While so-called *guided* data mining methods have been produced for a number of data mining areas including clustering⁶, association mining^{4,7}, and classification⁸, there is an architectural aspect to guided data mining, and to VDM in general, that has not been adequately explored and which represents an area for future work.

Another area of future work for the VDM community is quantification. Although the benefits that VDM can provide are clear to us, due to its subjective nature, the benefits of this synergy are not easily quantified and thus may not be as obvious to others. VDM methods can be more time-consuming to develop and thus for VDM to be accepted more widely we must find methods of showing that VDM demonstrates either (or both of) a time improvement or a quality improvement over non-visual methods.

This book has been long awaited. The VDM community has come a long way in a short time. Due to its ability to merge the cognitive ability and contextual awareness of humans with the increasing computational power of data mining systems, VDM is undoubtedly not just a future trend but destined to be one of the main themes for data mining for many years to come.

April 2008

John F. Roddick

⁴ Ceglar, A., Roddick, J.F.: GAM - a guidance enabled association mining environment. *International Journal of Business Intelligence and Data Mining* 2, 3–28 (2007)

⁵ Ceglar, A.: *Guided Association Mining through Dynamic Constraint Refinement*. PhD thesis, Flinders University (2005)

⁶ Anderson, D., Anderson, E., Lesh, N., Marks, J., Perlin, K., Ratajczak, D., Ryall, K.: Human guided simple search: combining information visualization and heuristic search. In: *Workshop on new paradigms in information visualization and manipulation; In conjunction with the 8th ACM international conference on Information and Knowledge Management, Kansas City, MO*, pp. 21–25. ACM Press, New York (2000)

⁷ Ng, R., Lakshmanan, L., Han, J., Pang, A.: Exploratory mining and pruning optimizations of constrained association rules. In: *17th ACM SIGACT-SIGMOD-SIGART Symposium on the Principles of Database Systems, Seattle, WA*, pp. 13–24. ACM Press, New York (1998)

⁸ Ankerst, M., Ester, M., Kriegel, H.P.: Towards an effective cooperation of the user and the computer for classification. In: *6th International Conference on Knowledge Discovery and Data Mining (KDD 2000), Boston, MA*, pp. 179–188 (2000)

Preface

John W. Tukey, who made unparalleled contributions to statistics and to science in general during his long career at Bell Labs and Princeton University, emphasized that seeing may be believing or disbelieving, but above all, data analysis involves visual, as well as statistical, understanding. Certainly one of the oldest visual explanations in mathematics is the visual proof of the Pythagorean theorem. The proof, impressive in its brevity and elegance, stresses the power of an interactive visual representation in facilitating our analytical thought processes. Thus, visual reasoning approaches to extracting and comprehension of the information encoded in data sets became the focus of what is called *visual data mining*. The field emerged from the integration of concepts from numerous fields, including computer graphics, visualization metaphors and methods, information and scientific data visualization, visual perception, cognitive psychology, diagrammatic reasoning, 3D virtual reality systems, multimedia and design computing, data mining and online analytical processing, very large databases last, and even collaborative virtual environments.

The importance of the field had already been recognized in the beginning of the decade. This was reflected in the series of visual data mining workshops, conducted at the major international conferences devoted to data mining. Later, the conferences and periodicals in information visualization paid substantial attention to some developments in the field. Commercial tools and the work in several advanced laboratories and research groups across the globe provided working environments for experimenting not only with different methods and techniques for facilitating the human visual system in examination and patterns discovery, and understanding of patterns among massive volumes of multi-dimensional and multi-source data, but also for testing techniques that provide robust and statistically valid visual patterns. It was not until a panel of more than two dozen internationally renowned individuals was assembled, in order to address the shortcomings and drawbacks of the current state of visual information processing, that the need for a systematic and methodological development of visual analytics was placed in the top priorities on the research and development agenda in 2005.

This book aims at addressing this need. Through a collection of 21 chapters selected from more than 46 submissions, it offers a systematic presentation of the state of the art in the field, presenting it in the context of visual analytics. Since visual analysis is such a different technique, it is an extremely significant topic for contemporary data mining and data analysis.

The editors would like to thank all the authors for their contribution to the volume and their patience in addressing reviewers' and editorial feedback. Without their contribution and support the creation of this volume would have been impossible. The editors would like to thank the reviewers for their thorough reviews and detailed comments.

Special thanks go to John Roddick, who, on short notice, kindly accepted the invitation to write the Foreword to the book.

April 2008

Simeon J. Simoff
Michael Böhlen
Arturas Mazeika

Table of Contents

Visual Data Mining: An Introduction and Overview	1
<i>Simeon J. Simoff, Michael H. Böhlen, and Arturas Mazeika</i>	
Part 1 – Theory and Methodologies	
The 3DVDM Approach: A Case Study with Clickstream Data	13
<i>Michael H. Böhlen, Linas Bukauskas, Arturas Mazeika, and Peer Mylov</i>	
Form-Semantics-Function – A Framework for Designing Visual Data Representations for Visual Data Mining	30
<i>Simeon J. Simoff</i>	
A Methodology for Exploring Association Models	46
<i>Alipio Jorge, João Poças, and Paulo J. Azevedo</i>	
Visual Exploration of Frequent Itemsets and Association Rules	60
<i>Li Yang</i>	
Visual Analytics: Scope and Challenges	76
<i>Daniel A. Keim, Florian Mansmann, Jörn Schneidewind, Jim Thomas, and Hartmut Ziegler</i>	
Part 2 – Techniques	
Using Nested Surfaces for Visual Detection of Structures in Databases	91
<i>Arturas Mazeika, Michael H. Böhlen, and Peer Mylov</i>	
Visual Mining of Association Rules	103
<i>Dario Bruzzese and Cristina Davino</i>	
Interactive Decision Tree Construction for Interval and Taxonomical Data	123
<i>François Poulet and Thanh-Nghi Do</i>	
Visual Methods for Examining SVM Classifiers	136
<i>Doina Caragea, Dianne Cook, Hadley Wickham, and Vasant Honavar</i>	
Text Visualization for Visual Text Analytics	154
<i>John Risch, Anne Kao, Stephen R. Poteet, and Y.-J. Jason Wu</i>	
Visual Discovery of Network Patterns of Interaction between Attributes	172
<i>Simeon J. Simoff and John Galloway</i>	

Mining Patterns for Visual Interpretation in a Multiple-Views Environment	196
<i>José F. Rodrigues Jr., Agma J.M. Traina, and Caetano Traina Jr.</i>	
Using 2D Hierarchical Heavy Hitters to Investigate Binary Relationships	215
<i>Daniel Trivellato, Arturas Mazeika, and Michael H. Böhlen</i>	
Complementing Visual Data Mining with the Sound Dimension: Sonification of Time Dependent Data	236
<i>Monique Noirhomme-Fraiture, Olivier Schöller, Christophe Demoulin, and Simeon J. Simoff</i>	
Context Visualization for Visual Data Mining	248
<i>Mao Lin Huang and Quang Vinh Nguyen</i>	
Assisting Human Cognition in Visual Data Mining	264
<i>Simeon J. Simoff, Michael H. Böhlen, and Arturas Mazeika</i>	
 Part 3 – Tools and Applications	
Immersive Visual Data Mining: The 3DVDM Approach	281
<i>Henrik R. Nagel, Erik Granum, Søren Bovbjerg, and Michael Vittrup</i>	
DataJewel: Integrating Visualization with Temporal Data Mining	312
<i>Mihael Ankerst, Anne Kao, Rodney Tjoelker, and Changzhou Wang</i>	
A Visual Data Mining Environment	331
<i>Stephen Kimani, Tiziana Catarci, and Giuseppe Santucci</i>	
Integrative Visual Data Mining of Biomedical Data: Investigating Cases in Chronic Fatigue Syndrome and Acute Lymphoblastic Leukaemia.	367
<i>Paul J. Kennedy, Simeon J. Simoff, Daniel R. Catchpoole, David B. Skillicorn, Franco Ubaudi, and Ahmad Al-Oqaily</i>	
Towards Effective Visual Data Mining with Cooperative Approaches	389
<i>François Poulet</i>	
 Author Index	 407