

**PROCEEDINGS OF
THE 2015 INTERNATIONAL CONFERENCE ON
DATA MINING**

DMIN 2015

Editors

**Robert Stahlbock
Gary M. Weiss**

Associate Editors

**Mahmoud Abou-Nasr
Hamid R. Arabnia**



WORLDCOMP'15

July 27-30, 2015

Las Vegas Nevada, USA

www.world-academy-of-science.org

©CSREA Press

This volume contains papers presented at The 2015 International Conference on Data Mining (DMIN'15). Their inclusion in this publication does not necessarily constitute endorsements by editors or by the publisher.

Copyright and Reprint Permission

Copying without a fee is permitted provided that the copies are not made or distributed for direct commercial advantage, and credit to source is given. Abstracting is permitted with credit to the source. Please contact the publisher for other copying, reprint, or republication permission.

Copyright © 2015 CSREA Press
ISBN: 1-60132-403-0
Printed in the United States of America

CSREA Press
U. S. A.

Foreword

It gives us great pleasure to introduce this collection of papers to be presented at the 11th International Conference on Data Mining 2015, DMIN'15 (www.dmin-2015.com), July 27-30, 2015, at Monte Carlo Resort, Las Vegas, Nevada, USA.

Data mining is a relatively young discipline that is critically important if we want to effectively learn from the tremendous amounts of data that are routinely being generated in science, engineering, medicine, business, and other areas in order to gain insight into processes, transactions, make better decisions, and deliver value to users or organizations. During the last years, we all observe new, more glorious and promising concepts or labels emerging and slowly but steadily displacing 'data mining' from the agenda of CTO's. It is the time of big data, advanced-/business-/customer-/data-/.../risk-analytics, to name only a few terms that dominate websites, trade journals, and the general press. But they all aim at leveraging data for a better understanding of and insight into complex real-world phenomena. They all pursue this objective using some formal, often algorithmic, procedures, at least to some extent. This is what data miners have been doing for decades. So maybe the label 'data mining' has lost much of its momentum and made room for more recent 'competitors', but the very idea of it, the idea to think of massive, omnipresent amounts of data as strategic assets, and the aim to capitalize on these assets by means of analytic procedures is, indeed, more relevant and topical than ever before. Advances in hardware and software are helpful, but there are still many challenges to be tackled in order to leverage the promises of data analytics.

An important mission of the World Congress in Computer Science, Computer Engineering, and Applied Computing (WORLDCOMP, a federated congress to which this conference is affiliated with) includes *"Providing a unique platform for a diverse community of constituents composed of scholars, researchers, developers, educators, and practitioners. The Congress makes concerted effort to reach out to participants affiliated with diverse entities (such as: universities, institutions, corporations, government agencies, and research centers/labs) from all over the world. The congress also attempts to connect participants from institutions that have teaching as their main mission with those who are affiliated with institutions that have research as their main mission."* By any definition of diversity, this congress is among the most diverse scientific meeting in USA. We are proud to report that this federated congress has authors and participants from 76 different nations representing variety of personal and scientific experiences that arise from differences in culture and values. The program committee of this conference as well as the program committee of all other tracks of the federated congress are as diverse as its authors and participants.

Data mining attracts innovative and influential contributions to both research and practice, across a wide range of academic disciplines and application domains. DMIN conferences seek to acknowledge and facilitate excellence in research and applications in the area of data mining. DMIN conferences are held annually within WORLDCOMP. WORLDCOMP'15 assembles a spectrum of 20 affiliated research conferences, workshops, and symposiums into a coordinated research meeting. Each conference has its own program committee as well as referees and own indexed proceedings. Attendees have full access to all 20 conferences' sessions, tracks, and tutorials. DMIN seeks to reflect the multi- and interdisciplinary nature of data mining and to facilitate the exchange and development of novel ideas, open communication and networking amongst researchers and practitioners in different research domains. As in previous years, we hope that the 2015 International Conference on Data Mining will provide a forum for you to present your research in a professional environment, exchange ideas, and network and interact across research areas. DMIN actively supports students and beginning researchers from lesser developed countries by funding registration and accommodation, in order to allow for a truly international networking and understanding. DMIN'15 provides an international and multicultural

experience with contributions from 20 different countries. We consider the resulting diversity in attendees and the mixture of established and starting researchers as a particular advantage of an engaging conference format.

DMIN'15 attracted a high number of submissions of theoretical research papers as well as industrial reports, application case studies, and in a second phase, late breaking papers, position papers, and abstract papers. The program committee would like to thank all those who submitted papers for consideration. We strived to establish a review process of high quality. To ensure a fair, objective and transparent review process all review criteria were published on the website. Papers were evaluated regarding their relevance to DMIN, originality, significance, information content, clarity, and soundness on an international level. Each aspect was objectively evaluated, with alternative aspects finding consideration for application papers. Each paper was refereed by at least two researchers in the topical area, taking the reviewers' expertise and confidence into consideration, with most of the papers receiving three reviews. The review process was competitive. The overall paper acceptance rate for papers was 49%.

We are very grateful to the many colleagues who helped in organizing the conference. In particular, we would like to thank the members of the program committee of DMIN'15 and the members of the congress steering committee. The continuing support of the DMIN program committee has been essential to further improve the quality of accepted submissions and the resulting success of the conference. The DMIN'15 program committee members are (in alphabetical order): Mahmoud Abou-Nasr, Jérôme Azé, James Buckley, Paulo Cortez, Kevin Daimi, Qin Ding, António Dourado, Philippe Fournier-Viger, Diego Galar, Peter Geczy, Zahid Halim, Tzung-Pei Hong, Wei-Chiang Hong, Sebastian Klenk, Terje Kristensen, Philippe Lenca, Stefan Lessmann, Wen-Yang Lin, Tanja Magoc, José M. Merigo Lindahl, Gerald Schaefer, Zhang Sen, Sabrina Senatore, Victor Sheng, Vijendra Singh, Robert Stahlbock, Ryszard Tadeusiewicz, Nicole Vincent, Baoying Wang, Chamont Wang, Simon Wang, Gary M. Weiss, Zijiang Yang, Yu Zhang, and Shang-Ming Zhou.

We would also like to thank our publicity co-chairs Ashu M. G. Solo (Fellow of British Computer Society, Principal/R&D Engineer, Maverick Technologies America Inc.) for circulating information on the conference, as well as www.KDnuggets.com, a platform for analytics, data mining and data science resources, for listing DMIN'15.

Considering the increasing efforts of all towards the quality of the review process, the conference sessions and the social program of DMIN'15, we are confident that you will find the conference stimulating and rewarding. It is a particular pleasure to provide data mining oriented invited talks and tutorials presented by the following esteemed members of the data mining community: Diego Galar (Luleå University of Technology, Sweden), and Peter Geczy (AIST, Japan).

The DMIN'15 conference organizers are also thankful to a number of co-sponsors, without whom the conference would not have been possible. As Sponsors-at-large, partners, and/or organizers each of the followings (separated by semicolons) provided help for at least one track of the World Congress: Computer Science Research, Education, and Applications Press (CSREA); US Chapter of World Academy of Science (<http://www.world-academy-of-science.org/>); American Council on Science & Education & Federated Research Council (<http://www.americancse.org/>); HoIP, Health Without Boundaries, Healthcare over Internet Protocol, UK (<http://www.hoip.eu>); HoIP Telecom, UK (<http://www.hoip-telecom.co.uk>); and WABT, Human Health Medicine, UNESCO NGOs, Paris, France (<http://www.thewabt.com/>). In addition, a number of university faculty members and their staff (names appear on the cover of the set of proceedings), several publishers of computer science and computer engineering books and journals, chapters and/or task forces of computer science associations/organizations from 4 countries, and developers of high-performance machines and systems provided significant help in organizing the conference as well as providing some resources. We are grateful to them all.

We are also grateful for support by the Institute of Information Systems at Hamburg University, Germany (www.uni-hamburg.de/IWI) and the Business Intelligence Laboratory, B I³S lab, Hamburg, Germany (www.bis-lab.com).

We express our gratitude to keynote, invited, and individual conference/tracks and tutorial speakers – the list of speakers appears on the conference web site. We would also like to thank the followings: UCMSS (Universal Conference Management Systems & Support, California, USA) for managing all aspects of the conference; Dr. Tim Field of APC for managing and coordinating the printing of the proceedings; and the staff of Monte Carlo Resort (Convention department) in Las Vegas for the professional service they provided. We would also like to thank Mahmoud Abou-Nasr for his continuous effort in organizing the special session on real-world data mining applications throughout the years.

Last but not least, we wish to express again our sincere gratitude and respect towards Prof. Hamid R. Arabnia (University of Georgia, USA); Coordinator of the federated congress, for his excellent and tireless support, organization and coordination of all affiliated events. His exemplary and professional effort in 2015 and all the years before in the WORLDCOMP steering committee makes these events possible!

Thank you all for your contribution to DMIN'15! We hope that you will experience a stimulating conference with many opportunities for future contacts, research and applications.

We present the proceedings of DMIN'15.

Robert Stahlbock
DMIN'15 General Conference Chair

Gary M. Weiss

Steering Committee DMIN'15
www.dmin-2015.com

Steering Committee WORLDCOMP, 2015
<http://www.world-academy-of-science.org/>

Contents

SESSION: REAL-WORLD DATA MINING APPLICATIONS, CHALLENGES, AND PERSPECTIVES

The Need for Big Data Collection and Analysis to Support the Development of an Advanced Maintenance Strategy 3

David Baglee, Salla Marttonen, Diego Galar

Financial Footnote Analysis: Developing a Text Mining Approach 10

Maryam Heidari, Carsten Felden

Product 's Quality Prediction with Respect to Equipments Data 17

Mariam Melhem, Bouchra Ananou, Mohand Djeziri, Mustapha Ouladsine, Jacques Pinaton

eMaintenance Platform for Performing Data Fusion Mutation on Machine Tools 24

Victor Simon, Diego Galar, David Baglee

Selecting a Classification Ensemble and Detecting Process Drift in an Evolving Data Stream 31

Alejandro Heredia-Langner, Luke Rodriguez, Andy Lin, Jennifer Webster

Reliability Evaluation of Underground Power Cables with Probabilistic Models 37

Hassan Mashad Nemati, Anita Sant'Anna, Sławomir Nowaczyk

Use of Social Networks Sites (SNSs) as A Collaborative Learning Technique: Survey Analysis and Mining Approach 44

Nevine Labib, Ahmed Sabry, Rasha Mostafa, Edward Morcos

Wavelet-Coupled Machine Learning Methods for Drought Forecast Utilizing Hybrid Meteorological and Remotely-Sensed Data 50

Robin Tan, Marek Perkowski

A Hierarchical Clustering Approach to Analyze Similarities between Sea Surface Temperature Patterns in the Caribbean 57

Marc Boumedine

Interactive Data Quality Assistance – An Approach for Min(d)ing the Quality of Data 62

Nadia El Bekri, Elisabeth Peinsipp-Byma

SESSION: SEGMENTATION, CLUSTERING, ASSOCIATION + WEB / TEXT / MULTIMEDIA MINING + SOFTWARE

AHC based Word Clustering considering Feature Similarity 67

Taeho Jo

Extraction of Relevant Entities in Textual Documents. Modeling Intelligence Maps	71
<i>Isnard Martins, Edgard Martins</i>	
Graph-based Link Prediction in Cross-session Task Identification	78
<i>Chao Xu, Mingzhu Zhu, Wei Xiong, Yi-fang Wu</i>	
Learning Temporal Regression Models and Voronoi Tessellation for Job Offers Recommendation	85
<i>Sidahmed Benabderrahmane, Nedra Mellouli, Myriam Lamolle, Jean-Baptiste Gabriel</i>	
Using Text Mining of Amazon Reviews to Explore User-defined Product Highlights and Issues	92
<i>Lleyana Jack, Yi-fang Tsai</i>	
How Can We Measure the Similarity Between Resumes of Selected Candidates for a Job?	99
<i>Luis Adrian Cabrera-Diego, Barthelemy Durette, Matthieu Lafon, Juan-Manuel Torres-Moreno, Marc El-Beze</i>	
A Language-Based and Process-Oriented Approach for Supporting the Knowledge Discovery Processes	107
<i>Hesham Mansour, Daniel Duchamp, Carl-Arndt Krapp</i>	
Exploiting Temporal Patterns of Hot Events in Weibo	116
<i>Jiakun Huang, Kai Niu, Zhiqiang He</i>	
SESSION: REGRESSION AND CLASSIFICATION	
Efficient Classifier over Stream Sliding Window using Associative Classification	125
<i>Prasanna Lakshmi Kompalli, Ramesh Kumar Reddy Cherku</i>	
Identifying Causes of Neonatal Mortality from Observational Data: A Bayesian Network Approach	132
<i>Kevin Wilson, Dennis Wallace, Shivaprasad Goudar, Douglas Theriaque, Elizabeth McClure</i>	
Learning Decision Trees from Histogram Data	139
<i>Ram Gurung, Tony Lindgren, Henrik Bostrom</i>	
Comparison Between Random Forest Algorithm and J48 Decision Trees Applied to the Classification of Power Quality Disturbances	146
<i>Fabbio Borges, Ricardo Fernandes, Lucas Moraes, Ivan Silva</i>	
Modelling Ground-Level Ozone Concentration using Ensemble Learning Algorithms	148
<i>Eman Al Abri, Eran Edirisinghe, Amin Nawadha</i>	

Automatic Detection Of Small Groups Of Persons, Influential Members, Relations And Hierarchy In Written Conversations Using Fuzzy Logic 155

French Pope, Rouzbeh Shirvani, Mugizi Rwebangira, Mohamed Chouikha, Ayo Taylor, Andres Ramirez, Amirsina Torfi

A Machine Learning Approach for Business Intelligence Analysis using Commercial Shipping Transaction Data 162

Lisa Bramer, Samrat Chatterjee, Aimee Holmes, Sean Robinson, Steven Bradley, Bobbie-Jo Webb-Robertson

SESSION: DATA MINING: LATE BREAKING PAPERS

Supervised Machine Learning Approach for Gender Disambiguation from Unstructured (Text) Documents 171

Amit Choudhary, Praveen Kumar, Sridhar Jeyaraman

Scalable Mining of Frequent and Significant Sequential Patterns 178

Zsolt T Kardkovacs, Gabor Kovacs

Raise Regression: Selection of the Raise Parameter 185

Catalina Garcia, Jose Garcia, Roman Salmeron, Maria Del Mar Lopez

