Lecture Notes in Artificial Intelligence 11828

Subseries of Lecture Notes in Computer Science

Series Editors

Randy Goebel University of Alberta, Edmonton, Canada Yuzuru Tanaka Hokkaido University, Sapporo, Japan Wolfgang Wahlster DFKI and Saarland University, Saarbrücken, Germany

Founding Editor

Jörg Siekmann DFKI and Saarland University, Saarbrücken, Germany More information about this series at http://www.springer.com/series/1244

Petra Kralj Novak · Tomislav Šmuc · Sašo Džeroski (Eds.)

Discovery Science

22nd International Conference, DS 2019 Split, Croatia, October 28–30, 2019 Proceedings



Editors Petra Kralj Novak Jožef Stefan Institute Ljubljana, Slovenia

Sašo Džeroski Jožef Stefan Institute Ljubljana, Slovenia Tomislav Šmuc Rudjer Bošković Institute Zagreb, Croatia

ISSN 0302-9743 ISSN 1611-3349 (electronic) Lecture Notes in Artificial Intelligence ISBN 978-3-030-33777-3 ISBN 978-3-030-33778-0 (eBook) https://doi.org/10.1007/978-3-030-33778-0

LNCS Sublibrary: SL7 - Artificial Intelligence

© Springer Nature Switzerland AG 2019

The chapter "Sparse Robust Regression for Explaining Classifiers" is Open Access. This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (http://creativecommons.org/licenses/by/4.0/). For further details see license information in the chapter.

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Preface

The Discovery Science conference presents a unique combination of latest advances in the development and analysis of methods for discovering scientific knowledge, coming from machine learning, data mining, and intelligent data analysis, with their application in various scientific domains.

The 22nd International Conference on Discovery Science (DS 2019) was held in Split, Croatia, during October 28–30, 2019. This was the first time the conference was organized as a stand-alone event. For its first 20 editions, DS was co-located with the International Conference on Algorithmic Learning Theory (ALT). In 2018 it was co-located with the 24th International Symposium on Methodologies for Intelligent Systems (ISMIS 2018).

DS 2019 received 63 international submissions. Each submission was reviewed by at least three Program Committee (PC) members. The PC decided to accept 21 regular papers and 19 short papers. This resulted in an acceptance rate of 33% for regular papers.

The conference included three keynote talks. Marinka Žitnik (Stanford University) contributed a talk titled "Representation Learning as a New Approach to Biomedical Research," Guido Caldarelli (IMT Lucca and ECLT Venice) gave a presentation titled "The Structure of Financial Networks," and Dino Pedreschi (University of Pisa), contributed a talk titled "Data and Algorithmic Bias: Explaining the Network Effect in Opinion Dynamics and the Training Data Bias in Machine Learning." Abstracts of the invited talks with short biographies of the invited speakers are included in these proceedings.

Besides the presentation of regular and short papers in the main program, the conference offered two new sessions. The "PhD Symposium" gave an opportunity to PhD students at an early stage of their studies to participate in the conference by presenting the topics of and early results from their research and discuss their work and experiences with peers, senior researchers and leading experts working on similar problems. The session titled "Late Breaking Contributions" featured poster and spotlight presentations of very recent research results on topics related to Discovery Science.

We are grateful to Springer for their long-term support, which got even stronger this year. Springer publishes the conference proceedings, as well as a regular special issue of the Machine Learning journal on Discovery Science. The latter offers authors a chance of publishing in this prestigious journal significantly extended and reworked versions of their DS conference papers, while being open to all submissions on DS conference topics.

This year, Springer (LNCS and Machine Learning journal), supported the best student paper awards. For DS 2019, the awardees are Anton Björklund, Andreas Henelius, Emilia Oikarinen, Kimmo Kallonen and Kai Puolamäki (for the paper "Sparse Robust Regression for Explaining Classifiers") and Yannik Klein, Michael Rapp and Eneldo Loza Mencía (for the paper "Efficient Discovery of Expressive Multi-label Rules Using Relaxed Pruning.") We would like to thank the Best Paper Award committee composed of Dragan Gamberger and Toon Calders for their precious and timely evaluations.

On the program side, we would like to thank all the authors of submitted papers, the PC members and the additional reviewers for their efforts in evaluating the submitted papers, as well as the keynote speakers. On the organization side, we would like to thank all the members of the Organizing Committee: Tomislav Lipić, Ana Vidoš, Matija Piškorec and Ratko Mileta, for the smooth preparation and organization of all conference associated activities. We are also grateful to the people behind EasyChair for developing the conference organization system that proved to be an essential tool in the paper submission and evaluation process, as well as in the preparation of the Springer proceedings.

The DS 2019 conference was organized under the auspices of the Rudjer Bošković Institute in Zagreb. The event was also supported by the Project of the Croatian Center for Excellence in Data Science and Advanced Cooperative Systems. Significant support, especially through human resources, was also provided by the Jožef Stefan Institute from Ljubljana. Finally, we are indebted to all conference participants, who contributed to making this exciting event a worthwhile endeavor for all involved.

October 2019

Petra Kralj Novak Tomislav Šmuc Sašo Džeroski

Organization

General Chair

Sašo Džeroski	Jožef Stefan Institute, Slovenia
Program Committee	Chairs
Petra Kralj Novak Tomislav Šmuc	Jožef Stefan Institute, Slovenia Rudjer Bošković Institute, Croatia
PhD Symposium Cha	air
Tomislav Lipić	Rudjer Bošković Institute, Croatia
Proceedings Chair	
Matija Piškorec	Rudjer Bošković Institute, Croatia
Web and Social Med	ia Chairs
Ratko Mileta Matija Piškorec	Rudjer Bošković Institute, Croatia Rudjer Bošković Institute, Croatia
Local Arrangements	Chair
Ana Vidoš	Rudjer Bošković Institute, Croatia
Program Committee	
Annalisa Appice	University of Bari Aldo Moro, Italy
Martin Atzmueller	Tilburg University, The Netherlands
Viktor Bengs	Paderborn University, Germany
Concha Bielza Lozoya	Universidad Politécnica de Madrid, Spain
Albert Bifet	LTCI, Telecom ParisTech, France
Alberto Cano	Virginia Commonwealth University, USA
Michelangelo Ceci	University of Bari Aldo Moro, Italy
Bruno Cremilleux	University of Caen Normandy, France
Claudia d'Amato	University of Bari Aldo Moro, Italy
Nicola Di Mauro	University of Bari Aldo Moro, Italy
Ivica Dimitrovski	Ss. Cyril and Methodius University in Skopje, North Macedonia

Wouter Duivesteijn Eindhoven University of Technology, The Netherlands Lina Fahed IMT Atlantique, France University of Oslo, Norway Hadi Fanaee University of Bari Aldo Moro, Italy Nicola Fanizzi Stefano Ferilli University of Bari Aldo Moro, Italy Johannes Fürnkranz Technische Universität Darmstadt, Germany Birmingham City University, UK Mohamed Gaber University of Porto, Portugal João Gama Rudier Bošković Institute, Croatia Dragan Gamberger Makoto Haraguchi Hokkaido University, Japan Kouichi Hirata Kyushu Institute of Technology, Japan Jaakko Hollmén Aalto University, Finland Paderborn University, Germany Eyke Huellermeier University of Porto, Portugal Alípio Jorge Ryukoku University, Japan Masahiro Kimura Jožef Stefan Institute, Slovenia Dragi Kocev Stefan Kramer Johannes Gutenberg University Mainz, Germany Ilia State University, Georgia Vincenzo Lagani Pedro Larranaga University of Madrid, Spain Nada Lavrač Jožef Stefan Institute, Slovenia Jurica Levatić Institute for Research in Biomedicine, Spain Tomislav Lipić Rudjer Bošković Institute, Croatia Francesca Alessandra Lisi University of Bari Aldo Moro, Italy Ss. Cyril and Methodius University in Skopje, Gjorgji Madjarov North Macedonia Institute for High Performance Computing Giuseppe Manco and Networking, Italy University of Rijeka, Croatia Sanda Martinčić-Ipšić Elio Masciari Institute for High Performance Computing and Networking, Italy University of Pisa, Italy Anna Monreale Siegfried Nijssen Université Catholique de Louvain, Belgium Rita P. Ribeiro University of Porto, Portugal Jožef Stefan Institute, Slovenia Panče Panov University of Torino, Italy Ruggero G. Pensa Bernhard Pfahringer University of Waikato, New Zealand Gianvito Pio University of Bari Aldo Moro, Italy Pascal Poncelet LIRMM Montpellier, France French Research Institute for Digital Sciences, France Jan Ramon French Research Institute for Digital Sciences, France Chedy Raïssi Marko Robnik-Šikonja University of Ljubljana, Slovenia University of Shizuoka, Japan Kazumi Saito Marina Sokolova University of Ottawa and Institute for Big Data Analytics, Canada Poznan University of Technology, Poland Jerzy Stefanowski

Ljupčo Todorovski Luis Torgo Herna Viktor Albrecht Zimmermann Blaž Zupan University of Ljubljana, Slovenia Dalhousie University, Canada University of Ottawa, Canada Université Caen Normandie, France University of Ljubljana, Slovenia

Additional Reviewers

Ahmadi, Mohsen Barracchia, Emanuele Pio Cancela, Brais Chambers, Lorraine Fernandes, Sofia Ghomeshi, Hossein Guarascio, Massimo Koptelov, Maksim Kulikovskikh, Ilona Oliveira, Mariana Pasquadibisceglie, Vincenzo Pisani, Francesco S. Stepišnik, Tomaž Tabassum, Shazia Tornede, Tanja Wever, Marcel Zopf, Markus

Keynote Talks

The Structure of Financial Networks

Guido Caldarelli

IMT School for Advanced Studies, Lucca and European Centre for Living Technology, Venice

Abstract. Financial inter-linkages play an important role in the emergence of financial instabilities and the formulation of systemic risk can greatly benefit from a network approach. In this talk, we focus on the role of linkages along the two dimensions of contagion and liquidity, and we discuss some insights that have recently emerged from network models. With respect to the issue of the determination of the optimal architecture of the financial system, models suggest that regulators have to look at the interplay of network topology, capital requirements, and market liquidity. With respect to the issue of the determination of systemically important financial institutions, the findings indicate that both from the point of view of contagion and from the point of view of liquidity provision, there is more to systemic importance than just size. In particular for contagion, the position of institutions in the network matters and their impact can be computed through stress tests even when there are no defaults in the system.

We present an overview of the use of networks in Finance and Economics. We show how this approach enables us to address important questions as, for example, the stability of financial systems and the systemic risk associated with the functioning of the interbank market. For example with DebtRank, a novel measure of systemic impact inspired by feedback-centrality we are able to measure the nodes that become systemically important at the peak of the crisis. Moreover, a systemic default could have been triggered even by small dispersed shocks. The results suggest that the debate on too-big-to-fail institutions should include the even more serious issue of too-central-to-fail. All these results are new in the field and allow for a better understanding and modelling of different Financial systems.

Keywords: Financial networks · Systemic risk · Interbank market

Short Biography of the Lecturer: Guido Caldarelli is Full Professor in Theoretical Physics at IMT School for Advanced Studies Lucca, and is Research associate at the European Centre for Living Technology, Venice. His main scientific activity is the study of networks, mostly analysis and modelling of financial networks. Author of more than 200 publication on the subject and three books, he is currently the president of the Complex Systems Society. He has been coordinator of the FET IP Project MULTIPLEX: Foundational Research on Multilevel Complex Networks and Systems (2012–2016), the FET OPEN Project FoC: Forecasting Financial Crises (2010–2014), and the FET OPEN Project COSIN: Coevolution and Self Organization in Complex

Networks (2002–2005). Guido Caldarelli received his Ph.D. from SISSA, after which he was a postdoc in the Department of Physics and School of Biology, University of Manchester. He then worked at the Theory of Condensed Matter Group, University of Cambridge. He returned to Italy as a lecturer at National Institute for Condensed Matter (INFM) and later as Primo Ricercatore in the Institute of Complex Systems of the National Research Council of Italy. In this period, he was also the coordinator of the Networks subproject, part of the Complexity Project, for the Fermi Centre. He also spent some terms at University of Fribourg (Switzerland) and in 2006 he has been visiting professor at École Normale Supérieure in Paris. More information and a complete CV are available at: http://www.guidocaldarelli.com.

Data and Algorithmic Bias: Explaining the Network Effect in Opinion Dynamics and the Training Data Bias in Machine Learning

Dino Pedreschi

Università di Pisa, Istituto di Scienza e Tecnologie dell'Informazione, Consiglio Nazionale delle Ricerche http://kdd.isti.cnr.it

Abstract. Data science and network science are creating novel means to study the complexity of our societies and to measure, understand and predict social phenomena. My talk gives an overview of recent research at the Knowledge Discovery (KDD) Lab in Pisa within the SoBigData.eu research infrastructure, targeted at explaining the effects of data and algorithmic bias in different domains, using both data-driven and model-driven arguments. First, I introduce a model showing how algorithmic bias instilled in an opinion diffusion process artificially vields increased polarisation, fragmentation and instability in a population. Second, I focus on the urgent open challenge of how to construct meaningful explanations of opaque AI/ML black-box decision systems, introducing the local-to-global framework for the explanation of ML classifiers as a way towards explainable AI. The two cases show how the combination of data-driven and model-driven interdisciplinary research has a huge potential to shed new light on complex phenomena like discrimination and polarisation, as well as to explain how decision making black-boxes, both human and artificial, actually work. I conclude with an account of the open data science paradigm pursued in SoBigData.eu Research Infrastructure and its importance for interdisciplinary data driven science that impacts societal challenges.

Keywords: Explainable AI · Data bias · Algorithmic bias

Short Biography of the Lecturer: Dino Pedreschi is a professor of computer science at the University of Pisa, and a pioneering scientist in data science. He co-leads the Pisa KDD Lab – Knowledge Discovery and Data Mining Laboratory http://kdd.isti.cnr. it, a joint research initiative of the University of Pisa and the Information Science and Technology Institute of the Italian National Research Council. His research focus is on big data analytics and mining and their impact on society. He is a founder of the Business Informatics MSc program at University of Pisa, a course targeted at the education of interdisciplinary data scientists, and of SoBigData.eu, the European H2020 Research Infrastructure "Big Data Analytics and Social Mining Ecosystem" www.sobigdata.eu. Dino has been a visiting scientist at Barabasi Lab (Center for

Complex Network Research) of Northeastern University, Boston, and earlier at the University of Texas at Austin, at CWI Amsterdam and at UCLA. In 2009, Dino received a Google Research Award for his research on privacy-preserving data mining. Dino is a member of the expert group in AI of the Italian Ministry of research and the director of the Data Science PhD program at Scuola Normale Superiore in Pisa. Dino is a co-PI of the 2019 ERC grant XAI – Science and technology for the explanation of AI decision making (PI: Fosca Giannotti).

Representation Learning as a New Approach to Biomedical Research

Marinka Žitnik

Computer Science Department, School of Engineering, Stanford University

Abstract. Large datasets are being generated that can transform science and medicine. New machine learning methods are necessary to unlock these data and open doors for scientific discoveries. In this talk, I will argue that machine learning models should not be trained in the context of one particular dataset. Instead, we should be developing methods that combine data in their broadest sense into knowledge networks, enhance these networks to reduce biases and uncertainty, and then learn and reason over the networks. My talk will focus on two key aspects of this goal: representation learning and network science for knowledge networks. I will show how realizing this goal can set sights on new frontiers beyond classic applications of neural networks on biomedical image and sequence data. I will start by presenting a framework that learns deep models by embedding knowledge networks into compact embedding spaces whose geometry is optimized to reflect network topology, the essence of networks. I will then describe two applications of the framework to drug discovery and medicine. First, the framework allowed us to, for the first time, predict the safety of drug combinations at scale. We embedded a knowledge network of molecular, drug, and patient data at the scale of billions of interactions for all medications in the U.S. Using the embeddings, the approach can predict unwanted side effects for any combination of drugs that patients take, and we can validate predictions in the clinic using real patient data. Second, I will discuss how the framework enabled us to predict what diseases a new drug could treat. I will show how the new approach can make correct predictions for many recently repurposed drugs and can operate even on the hardest, yet critical, diseases for which no good treatments exist. I will conclude with future directions for learning over interaction data and translation of machine learning methods into solutions for biomedical problems.

Keywords: Biomedicine · Representation learning · Network science · Knowledge graphs

Short Biography of the Lecturer: Marinka Žitnik is a postdoctoral scholar in Computer Science at Stanford University. She will join Harvard University as a tenure-track assistant professor in December 2019. Her research investigates machine learning for sciences. Her methods have had a tangible impact in biology, genomics, and drug discovery, and are used by major biomedical institutions, including Baylor College of Medicine, Karolinska Institute, Stanford Medical School, and Massachusetts General Hospital. She received her Ph.D. in Computer Science from University of

xviii M. Žitnik

Ljubljana while also researching at Imperial College London, University of Toronto, Baylor College of Medicine, and Stanford University. Her work received several best paper, poster, and research awards from the International Society for Computational Biology. She was named a Rising Star in EECS by MIT and also a Next Generation in Biomedicine by The Broad Institute of Harvard and MIT, being the only young scientist who received such recognition in both EECS and Biomedicine. She is also a member of the Chan Zuckerberg Biohub at Stanford.

Contents

Advanced Machine Learning

The CURE for Class Imbalance	3
Mining a Maximum Weighted Set of Disjoint Submatrices Vincent Branders, Guillaume Derval, Pierre Schaus, and Pierre Dupont	18
Dataset Morphing to Analyze the Performance of Collaborative Filtering André Correia, Carlos Soares, and Alípio Jorge	29
Construction of Histogram with Variable Bin-Width Based on Change Point Detection	40
A Unified Approach to Biclustering Based on Formal Concept Analysis and Interval Pattern Structure	51
A Sampling-Based Approach for Discovering Subspace Clusters Sandy Moens, Boris Cule, and Bart Goethals	61
Epistemic Uncertainty Sampling Vu-Linh Nguyen, Sébastien Destercke, and Eyke Hüllermeier	72
Utilizing Hierarchies in Tree-Based Online Structured Output Prediction Aljaž Osojnik, Panče Panov, and Sašo Džeroski	87
On the Trade-Off Between Consistency and Coverage in Multi-label Rule Learning Heuristics	96
Hyperparameter Importance for Image Classification by Residual Neural Networks	112

Applications

Cellular Traffic Prediction and Classification: A Comparative Evaluation	
of LSTM and ARIMA	129
Amin Azari, Panagiotis Papapetrou, Stojan Denic, and Gunnar Peters	

XX	Contents
vv	Contents
AA	contents

Main Factors Driving the Open Rate of Email Marketing Campaigns Andreia Conceição and João Gama	145
Enhancing BMI-Based Student Clustering by Considering Fitness as Key Attribute Erik Dovgan, Bojan Leskošek, Gregor Jurak, Gregor Starc, Maroje Sorić, and Mitja Luštrek	155
Deep Learning Does Not Generalize Well to Recognizing Cats and Dogs in Chinese Paintings <i>Qianqian Gu and Ross King</i>	166
Temporal Analysis of Adverse Weather Conditions Affecting Wheat Production in Finland	176
Predicting Thermal Power Consumption of the Mars Express Satellite with Data Stream Mining Bozhidar Stevanoski, Dragi Kocev, Aljaž Osojnik, Ivica Dimitrovski, and Sašo Džeroski	186
Data and Knowledge Representation	
Parameter-Less Tensor Co-clustering Elena Battaglia and Ruggero G. Pensa	205
Deep Triplet-Driven Semi-supervised Embedding Clustering Dino Ienco and Ruggero G. Pensa	220
Neurodegenerative Disease Data Ontology Ana Kostovska, Ilin Tolovski, Fatima Maikore, the Alzheimer's Disease Neuroimaging Initiative, Larisa Soldatova, and Panče Panov	235
Embedding to Reference t-SNE Space Addresses Batch Effectsin Single-Cell ClassificationPavlin G. Poličar, Martin Stražar, and Blaž Zupan	246
Symbolic Graph Embedding Using Frequent Pattern Mining Blaž Škrlj, Nada Lavrač, and Jan Kralj	261
Feature Importance	

Feature Selection for Analogy-Based Learning to Rank	279
Mohsen Ahmadi Fahandar and Eyke Hüllermeier	

Contents xxi

Ensemble-Based Feature Ranking for Semi-supervised Classification	290
Matej Petković, Sašo Džeroski, and Dragi Kocev	

Variance-Based Feature Importance in Neural Networks	306
Cláudio Rebelo de Sá	

Interpretable Machine Learning

A Density Estimation Approach for Detecting and Explaining Exceptional Values in Categorical Data Fabrizio Angiulli, Fabio Fassetti, Luigi Palopoli, and Cristina Serrao	319
A Framework for Human-Centered Exploration of Complex Event Log Graphs	335
Sparse Robust Regression for Explaining Classifiers Anton Björklund, Andreas Henelius, Emilia Oikarinen, Kimmo Kallonen, and Kai Puolamäki	351
Efficient Discovery of Expressive Multi-label Rules Using Relaxed Pruning Yannik Klein, Michael Rapp, and Eneldo Loza Mencía	367

Networks

Evolving Social Networks Analysis via Tensor Decompositions: From Global Event Detection Towards Local Pattern Discovery and Specification	385
Efficient and Accurate Non-exhaustive Pattern-Based Change Detection in Dynamic Networks	396
A Combinatorial Multi-Armed Bandit Based Method for Dynamic Consensus Community Detection in Temporal Networks Domenico Mandaglio and Andrea Tagarelli	412
Resampling-Based Framework for Unbiased Estimator of Node Centrality over Large Complex Network	428

Pattern Discovery

Layered Learning for Early Anomaly Detection: Predicting Critical	
Health Episodes	445
Ensemble Clustering for Novelty Detection in Data Streams Kemilly Dearo Garcia, Elaine Ribeiro de Faria, Cláudio Rebelo de Sá, João Mendes-Moreira, Charu C. Aggarwal, André C. P. L. F. de Carvalho, and Joost N. Kok	460
Mining Patterns in Source Code Using Tree Mining Algorithms Hoang Son Pham, Siegfried Nijssen, Kim Mens, Dario Di Nucci, Tim Molderez, Coen De Roover, Johan Fabry, and Vadim Zaytsev	471
KnowBots: Discovering Relevant Patterns in Chatbot Dialogues Adriano Rivolli, Catarina Amaral, Luís Guardão, Cláudio Rebelo de Sá, and Carlos Soares	481
Fast Distance-Based Anomaly Detection in Images Using an Inception-Like Autoencoder Natasa Sarafijanovic-Djukic and Jesse Davis	493
Time Series	
Adaptive Long-Term Ensemble Learning from Multiple High-Dimensional Time-Series Samaneh Khoshrou and Mykola Pechenizkiy	511
Fourier-Based Parametrization of Convolutional Neural Networks for Robust Time Series Forecasting Sascha Krstanovic and Heiko Paulheim	522
Integrating LSTMs with Online Density Estimation for the Probabilistic Forecast of Energy Consumption	533