# Intelligent Systems Reference Library

Volume 149

The aim of this series is to publish a Reference Library, including novel advances and developments in all aspects of Intelligent Systems in an easily accessible and well structured form. The series includes reference works, handbooks, compendia, textbooks, well-structured monographs, dictionaries, and encyclopedias. It contains well integrated knowledge and current information in the field of Intelligent Systems. The series covers the theory, applications, and design methods of Intelligent Systems. Virtually all disciplines such as engineering, computer science, avionics, business, e-commerce, environment, healthcare, physics and life science are included. The list of topics spans all the areas of modern intelligent systems such as: Ambient intelligence, Computational intelligence, Social intelligence, Computational neuroscience, Artificial life, Virtual society, Cognitive systems, DNA and immunity-based systems, e-Learning and teaching, Human-centred computing and Machine ethics, Intelligent control, Intelligent data analysis, Knowledge-based paradigms, Knowledge management, Intelligent agents, Intelligent decision making, Intelligent network security, Interactive entertainment, Learning paradigms, Recommender systems, Robotics and Mechatronics including human-machine teaming, Self-organizing and adaptive systems, Soft computing including Neural systems, Fuzzy systems, Evolutionary computing and the Fusion of these paradigms, Perception and Vision, Web intelligence and Multimedia.

More information about this series at http://www.springer.com/series/8578

George A. Tsihrintzis · Dionisios N. Sotiropoulos
Lakhmi C. Jain

Editors

# Machine Learning Paradigms

Advances in Data Analytics

Springer

*Editors*
George A. Tsihrintzis
University of Piraeus
Piraeus
Greece

Dionisios N. Sotiropoulos
University of Piraeus
Piraeus
Greece

Lakhmi C. Jain
Faculty of Engineering and Information
    Technology, Centre for Artificial
    Intelligence
University of Technology
Sydney, NSW
Australia

and

Faculty of Science, Technology
    and Mathematics
University of Canberra
Canberra, ACT
Australia

and

KES International
Shoreham-by-Sea
UK

*To my wife and colleague, Prof.-Dr. Maria Virvou, and to our daughters, Evina, Konstantina and Andreani*

George A. Tsihrintzis

*To my beloved family and friends*

Dionisios N. Sotiropoulos

*To my beloved family*

Lakhmi C. Jain

# Foreword

In 1959, Arthur Samuel (1901–1990) published *Some Studies in Machine Learning Using the Game of Checkers* [1]. The paper was one of the earliest uses of the words "machine learning" [2]. He wrote, "As a result of these experiments one can say with some certainty that it is now possible to devise learning schemes which will greatly outperform an average person and that such learning schemes may eventually be economically feasible as applied to real-life problems" [1, p. 548]. His program together with IBM's first stored-program computer, the 701, demonstrated this statement by winning a game of checkers against a human expert in Connecticut. Since that time, games have provided fertile research ground for artificial intelligence, and in 1996, the Chinook Project for checkers was recognized by the Guinness Book of World Records as the first computer program to win a human world championship [3]. The day to apply machine learning to challenging real-world problems is here and now.

What is "machine learning"? Several suggested definitions are discussed on the IBM community site [4], including "The purpose of machine learning is to learn from training data in order to make as good as possible predictions on new, unseen, data." This definition suggests some of the challenges with machine learning. The program needs to build a model based on training data that includes the correct answer (i.e., supervised learning) and that minimizes the error in predicting new data. Alternatively, algorithms may look for structure in the data and group similar clusters (i.e., unsupervised learning). Too closely mirroring the training data results in overfitting and poor results with unknown data, and too little fitting results in unacceptable errors in predictions. In addition, as updated known data become available, the model may need to be re-adjusted to retain generalizability. Thus, the methods used in machine learning are constantly being researched and assessed against real-life data from various fields along with the computer technologies needed to implement them.

This book applies and assesses machine learning for classes of important real-life problems, an area often referred to as "data analytics." Authors have contributed leading research in the areas of medical, biological and signal sciences; social studies and social interactions; traffic, computer and power networks; and

digital forensics. The book also looks to the future for research areas that may yield theoretical advances. The editors have provided a valuable and much-needed collection of leading research in machine learning and data analytics that will increasingly impact each of us in our everyday lives. New and experienced researchers, practitioners, and those interested in machine learning will be inspired by the innovative ideas contained in its pages.

Baltimore, USA                                              Gloria Phillips-Wren, Ph.D.
                                              Professor, Loyola University Maryland

## References

1. Samuel, A.: Some studies in machine learning using the game of checkers. IBM J. 3(3), 210–229 (1959)
2. McCarthy, J., Feigenbaum, E.: In memoriam—Arthur Samuel: Pioneer in machine learning. AI Mag. 11(3), 10–11 (1990)
3. Chinook,: Arthur Samuel's legacy. Accessed on 28 May 2018, from https://webdocs.cs.ualberta.ca/∼chinook/project/legacy.html (2018)
4. Puget, J-C.: What is machine learning? IBM Community, May 18 from https://www.ibm.com/developerworks/community/blogs/jfp/entry/What_Is_Machine_Learning?lang=en (2016). Accessed on 28 May 2018

# Preface

At the dawn of the fourth Industrial Revolution, *data analytics* is emerging as a force that drives towards dramatic changes in our daily lives, the workplace and human relations. Synergies between physical, digital, biological and energy sciences and technologies, sewn together by *non-traditional data collection and analysis*, drive the digital economy at all levels and offer new, previously unavailable opportunities.

The need for data analytics arises in most modern scientific disciplines, including engineering, natural, computer and information sciences, economics, business, commerce, environment, healthcare and life sciences. The book at hand explores some of the emerging scientific and technological areas in which data analytics arises as a need and, thus, may play a significant role in the years to come.

Coming as the third volume under the general title *Machine Learning Paradigms* and following two related monographs, the book includes an editorial note (Chap. 1) and an additional twelve (12) chapters and is divided into five parts, namely: (1) *Data Analytics in the Medical, Biological and Signal Sciences*, (2) *Data Analytics in Social Studies and Social Interactions*, (3) *Data Analytics in Traffic, Computer and Power Networks*, (4) *Data Analytics for Digital Forensics* and (5) *Theoretical Advances and Tools for Data Analytics*.

This research book is directed towards professors, researchers, scientists, engineers and students of all disciplines. We hope that they all will find it useful in their works and researches.

We are grateful to the authors and the reviewers for their excellent contributions and visionary ideas. We are also thankful to Springer for agreeing to publish this book. Last, but not least, we are grateful to the Springer staff for their excellent work in producing this book.

Piraeus, Greece                                                    George A. Tsihrintzis
Piraeus, Greece                                              Dionisios N. Sotiropoulos
Sydney, Australia/Canberra, Australia                                 Lakhmi C. Jain

# Contents

**5 Machine Learning Methods for the Protein Fold Recognition Problem** ........................................... 101
Katarzyna Stapor, Irena Roterman-Konieczna and Piotr Fabian

**6 Speech Analytics Based on Machine Learning** ............... 129
Grazina Korvel, Adam Kurowski, Bozena Kostek
and Andrzej Czyzewski

## Part II  Data Analytics in Social Studies and Social Interactions

## Part III  Data Analytics in Traffic, Computer and Power Networks

**Part IV   Data Analytics for Digital Forensics**

**12  Combining Genetic Algorithms and Neural Networks for File
    Forgery Detection** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .  317
    Konstantinos Karampidis, Ioannis Deligiannis
    and Giorgos Papadourakis

**Part V   Theoretical Advances and Tools for Data Analytics**

**13  Deep Learning Analytics** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .  339
    Nikolaos Passalis and Anastasios Tefas