

Intelligent Systems Reference Library, Volume 12

Editors-in-Chief

Prof. Janusz Kacprzyk
Systems Research Institute
Polish Academy of Sciences
ul. Newelska 6
01-447 Warsaw
Poland
E-mail: kacprzyk@ibspan.waw.pl

Prof. Lakhmi C. Jain
University of South Australia
Adelaide
Mawson Lakes Campus
South Australia 5095
Australia
E-mail: Lakhmi.jain@unisa.edu.au

Further volumes of this series can be found on our
homepage: springer.com

Vol. 1. Christine L. Mumford and Lakhmi C. Jain (Eds.)
*Computational Intelligence: Collaboration, Fusion
and Emergence*, 2009
ISBN 978-3-642-01798-8

Vol. 2. Yuehui Chen and Ajith Abraham
*Tree-Structure Based Hybrid
Computational Intelligence*, 2009
ISBN 978-3-642-04738-1

Vol. 3. Anthony Finn and Steve Scheduling
*Developments and Challenges for
Autonomous Unmanned Vehicles*, 2010
ISBN 978-3-642-10703-0

Vol. 4. Lakhmi C. Jain and Chee Peng Lim (Eds.)
*Handbook on Decision Making: Techniques
and Applications*, 2010
ISBN 978-3-642-13638-2

Vol. 5. George A. Anastassiou
Intelligent Mathematics: Computational Analysis, 2010
ISBN 978-3-642-17097-3

Vol. 6. Ludmila Dymowa
Soft Computing in Economics and Finance, 2011
ISBN 978-3-642-17718-7

Vol. 7. Gerasimos G. Rigatos
Modelling and Control for Intelligent Industrial Systems, 2011
ISBN 978-3-642-17874-0

Vol. 8. Edward H.Y. Lim, James N.K. Liu, and
Raymond S.T. Lee
*Knowledge Seeker – Ontology Modelling for Information
Search and Management*, 2011
ISBN 978-3-642-17915-0

Vol. 9. Menahem Friedman and Abraham Kandel
Calculus Light, 2011
ISBN 978-3-642-17847-4

Vol. 10. Andreas Tolk and Lakhmi C. Jain
Intelligence-Based Systems Engineering, 2011
ISBN 978-3-642-17930-3

Vol. 11. Samuli Niiranen and Andre Ribeiro (Eds.)
Information Processing and Biological Systems, 2011
ISBN 978-3-642-19620-1

Vol. 12. Florin Gorunescu
Data Mining, 2011
ISBN 978-3-642-19720-8

Florin Gorunescu

Data Mining

Concepts, Models and Techniques

Prof. Florin Gorunescu
Chair of Mathematics
Biostatistics and Informatics University of
Medicine and Pharmacy of Craiova
Professor associated to the Department of
Computer Science
Faculty of Mathematics and Computer Science
University of Craiova
Romania
E-mail: gorun@umfcv.ro

ISBN 978-3-642-19720-8

e-ISBN 978-3-642-19721-5

DOI 10.1007/978-3-642-19721-5

Intelligent Systems Reference Library

ISSN 1868-4394

Library of Congress Control Number: 2011923211

© 2011 Springer-Verlag Berlin Heidelberg

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, reuse of illustrations, recitation, broadcasting, reproduction on microfilm 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.

The use of general descriptive names, registered names, trademarks, 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.

Typeset & Cover Design: Scientific Publishing Services Pvt. Ltd., Chennai, India.

Printed on acid-free paper

9 8 7 6 5 4 3 2 1

springer.com

To my family

Preface

Data Mining represents a complex of technologies that are rooted in many disciplines: mathematics, statistics, computer science, physics, engineering, biology, etc., and with diverse applications in a large variety of different domains: business, health care, science and engineering, etc. Basically, data mining can be seen as the science of exploring large datasets for extracting implicit, previously unknown and potentially useful information.

My aim in writing this book was to provide a friendly and comprehensive guide for those interested in exploring this vast and fascinating domain. Accordingly, my hope is that after reading this book, the reader will feel the need to deepen each chapter to learn more details.

This book aims to review the main techniques used in data mining, the material presented being supported with various examples, suggestively illustrating each method.

The book is aimed at those wishing to be initiated in data mining and to apply its techniques to practical applications. It is also intended to be used as an introductory text for advanced undergraduate-level or graduate-level courses in computer science, engineering, or other fields. In this regard, the book is intended to be largely self-contained, although it is assumed that the potential reader has a quite good knowledge of mathematics, statistics and computer science.

The book consists of six chapters, organized as follows:

- The first chapter introduces and explains fundamental aspects about data mining used throughout the book. These are related to: what is data mining, why to use data mining, how to mine data? Data mining solvable problems, issues concerning the modeling process and models, main data mining applications, methodology and terminology used in data mining are also discussed.
- Chapter 2 is dedicated to a short review regarding some important issues concerning data: definition of data, types of data, data quality, and types of data attributes.

- Chapter 3 deals with the problem of data analysis. Having in mind that data mining is an analytic process designed to explore large amounts of data in search of consistent and valuable hidden knowledge, the first step consists in an initial data exploration and data preparation. Then, depending on the nature of the problem to be solved, it can involve anything from simple descriptive statistics to regression models, time series, multivariate exploratory techniques, etc. The aim of this chapter is therefore to provide an overview of the main topics concerning exploratory data analysis.
- Chapter 4 presents a short overview concerning the main steps in building and applying classification and decision trees in real-life problems.
- Chapter 5 summarizes some well-known data mining techniques and models, such as: Bayesian and rule-based classifiers, artificial neural networks, k -nearest neighbors, rough sets, clustering algorithms, and genetic algorithms.
- The final chapter discusses the problem of evaluating the performance of different classification (and decision) models.

An extensive bibliography is included, which is intended to provide the reader with useful information covering all the topics approached in this book.

The organization of the book is fairly flexible, the selection of the topics to be approached being determined by the reader himself (herself), although my hope is that the book will be read entirely.

Finally, I wish this book to be considered just as a “compass” helping the interested reader to sail in the rough sea representing the current information vortex.

December 2010

Florin Gorunescu
Craiova

Contents

| | | |
|----------|--|----|
| 1 | Introduction to Data Mining | 1 |
| 1.1 | What Is and What Is Not Data Mining? | 1 |
| 1.2 | Why Data Mining? | 5 |
| 1.3 | How to Mine the Data? | 7 |
| 1.4 | Problems Solvable with Data Mining | 14 |
| 1.4.1 | Classification | 15 |
| 1.4.2 | Cluster Analysis | 19 |
| 1.4.3 | Association Rule Discovery | 23 |
| 1.4.4 | Sequential Pattern Discovery | 25 |
| 1.4.5 | Regression | 25 |
| 1.4.6 | Deviation/Anomaly Detection | 26 |
| 1.5 | About Modeling and Models | 26 |
| 1.6 | Data Mining Applications | 38 |
| 1.7 | Data Mining Terminology | 42 |
| 1.8 | Privacy Issues | 42 |
| 2 | The “Data-Mine” | 45 |
| 2.1 | What Are Data? | 45 |
| 2.2 | Types of Datasets | 46 |
| 2.3 | Data Quality | 50 |
| 2.4 | Types of Attributes | 52 |
| 3 | Exploratory Data Analysis | 57 |
| 3.1 | What Is Exploratory Data Analysis? | 57 |
| 3.2 | Descriptive Statistics | 59 |
| 3.2.1 | Descriptive Statistics Parameters | 60 |
| 3.2.2 | Descriptive Statistics of a Couple of Series | 68 |
| 3.2.3 | Graphical Representation of a Dataset | 81 |
| 3.3 | Analysis of Correlation Matrix | 85 |

| | | |
|----------|--|------------|
| 3.4 | Data Visualization | 89 |
| 3.5 | Examination of Distributions | 99 |
| 3.6 | Advanced Linear and Additive Models | 105 |
| 3.6.1 | Multiple Linear Regression | 105 |
| 3.6.2 | Logistic Regression | 116 |
| 3.6.3 | Cox Regression Model | 120 |
| 3.6.4 | Additive Models | 123 |
| 3.6.5 | Time Series: Forecasting | 124 |
| 3.7 | Multivariate Exploratory Techniques | 130 |
| 3.7.1 | Factor Analysis | 130 |
| 3.7.2 | Principal Components Analysis | 133 |
| 3.7.3 | Canonical Analysis | 136 |
| 3.7.4 | Discriminant Analysis | 137 |
| 3.8 | OLAP | 138 |
| 3.9 | Anomaly Detection | 148 |
| 4 | Classification and Decision Trees | 159 |
| 4.1 | What Is a Decision Tree? | 159 |
| 4.2 | Decision Tree Induction | 161 |
| 4.2.1 | GINI Index | 166 |
| 4.2.2 | Entropy | 169 |
| 4.2.3 | Misclassification Measure | 171 |
| 4.3 | Practical Issues Regarding Decision Trees | 179 |
| 4.3.1 | Predictive Accuracy | 179 |
| 4.3.2 | STOP Condition for Split | 179 |
| 4.3.3 | Pruning Decision Trees | 180 |
| 4.3.4 | Extracting Classification Rules from Decision Trees | 182 |
| 4.4 | Advantages of Decision Trees | 183 |
| 5 | Data Mining Techniques and Models | 185 |
| 5.1 | Data Mining Methods | 185 |
| 5.2 | Bayesian Classifier | 186 |
| 5.3 | Artificial Neural Networks | 191 |
| 5.3.1 | Perceptron | 192 |
| 5.3.2 | Types of Artificial Neural Networks | 205 |
| 5.3.3 | Probabilistic Neural Networks | 217 |
| 5.3.4 | Some Neural Networks Applications | 224 |
| 5.3.5 | Support Vector Machines | 234 |
| 5.4 | Association Rule Mining | 249 |
| 5.5 | Rule-Based Classification | 252 |
| 5.6 | k -Nearest Neighbor | 256 |
| 5.7 | Rough Sets | 260 |
| 5.8 | Clustering | 271 |
| 5.8.1 | Hierarchical Clustering | 282 |

| | | |
|--------------|---|------------|
| 5.8.2 | Non-hierarchical/Partitional Clustering | 284 |
| 5.9 | Genetic Algorithms | 289 |
| 5.9.1 | Components of GAs | 292 |
| 5.9.2 | Architecture of GAs | 310 |
| 5.9.3 | Applications | 313 |
| 6 | Classification Performance Evaluation | 319 |
| 6.1 | Costs and Classification Accuracy | 319 |
| 6.2 | ROC (Receiver Operating Characteristic) Curve | 323 |
| 6.3 | Statistical Methods for Comparing Classifiers | 328 |
| | References | 331 |
| Index | | 353 |