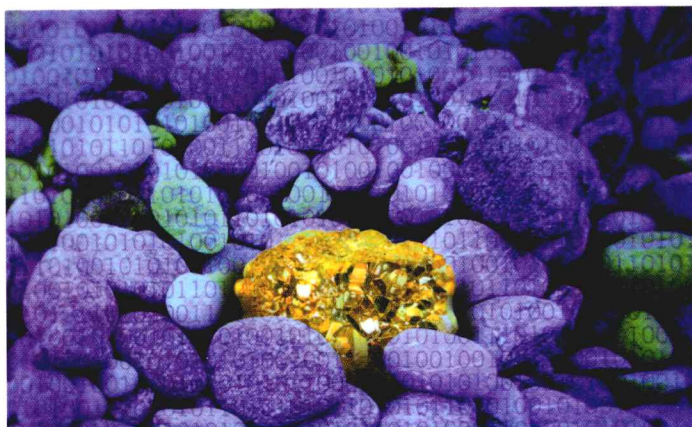


# 数据挖掘导论

(英文版)



INTRODUCTION TO **DATA MINING**



Pang-Ning Tan

密歇根州立大学

Michael Steinbach

明尼苏达大学

Vipin Kumar

明尼苏达大学

(美)

著



机械工业出版社  
China Machine Press

经典原版书库

# 数据挖掘导论

(英文版)

English reprint edition copyright © 2010 by Pearson Education Asia Limited and China Machine Press.

Original English language title: *Introduction to Data Mining* (ISBN 978-0-321-32136-7) by Pang-Ning Tan, Michael Steinbach and Vipin Kumar, Copyright © 2006.

All rights reserved.

Published by arrangement with the original publisher, Pearson Education, Inc., publishing as Addison-Wesley.

For sale and distribution in the People's Republic of China exclusively (except Taiwan, Hong Kong SAR and Macau SAR).

本书英文影印版由 Pearson Education Asia Ltd. 授权机械工业出版社独家出版。未经出版者书面许可, 不得以任何方式复制或抄袭本书内容。

仅限于中华人民共和国境内(不包括中国香港、澳门特别行政区和中国台湾地区)销售发行。

本书封面贴有 Pearson Education (培生教育出版集团) 激光防伪标签, 无标签者不得销售。

封底无防伪标均为盗版

版权所有, 侵权必究

本书法律顾问 北京市展达律师事务所

**本书版权登记号: 图字: 01-2010-4829**

**图书在版编目 (CIP) 数据**

数据挖掘导论 (英文版)/(美) 谭 (Tan, P. N.), 斯坦巴克 (Steinbach, M.), 库马尔 (Kumar, V.) 著. —北京: 机械工业出版社, 2010.9  
(经典原版书库)

书名原文: *Introduction to Data Mining*

ISBN 978-7-111-31670-1

I. 数… II. ①谭… ②斯… ③库… III. 数据采集 - 英文 IV. TP274

中国版本图书馆 CIP 数据核字 (2010) 第 170285 号

机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码 100037)

责任编辑: 李俊竹

北京京师印务有限公司印刷

2010 年 9 月第 1 版第 1 次印刷

150mm × 214mm · 24.625 印张

标准书号: ISBN 978-7-111-31670-1

定价: 59.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88378991; 88361066

购书热线: (010) 68326294; 88379649; 68995259

投稿热线: (010) 88379604

读者信箱: hzjsj@hzbook.com

## 出版者的话

文艺复兴以降，源远流长的科学精神和逐步形成的学术规范，使西方国家在自然科学的各个领域中取得了垄断性的优势；也正是这样的传统，使美国在信息技术发展的六十多年间名家辈出、独领风骚。在商业化的进程中，美国的产业界与教育界越来越紧密地结合，计算机学科中的许多泰山北斗同时身处科研和教学的最前线，由此而产生的经典科学著作，不仅擘划了研究的范畴，还揭示了学术的源变，既遵循学术规范，又自有学者个性，其价值并不会因年月的流逝而减退。

近年，在全球信息化大潮的推动下，我国的计算机产业发展迅猛，对专业人才的需求日益迫切。这对计算机教育界和出版界都既是机遇，也是挑战；而专业教材的建设在教育战略上显得举足轻重。在我国信息技术发展时间较短的现状下，美国等发达国家在其计算机科学发展的几十年间积淀和发展的经典教材仍有许多值得借鉴之处。因此，引进一批国外优秀计算机教材将对我国计算机教育事业的发展起到积极的推动作用，也是与世界接轨、建设真正的世界一流大学的必由之路。

机械工业出版社华章公司较早意识到“出版要为教育服务”。自1998年开始，我们就将工作重点放在了遴选、移译国外优秀教材上。经过多年的不懈努力，我们与 Pearson, McGraw-Hill, Elsevier, MIT, John Wiley & Sons, Cengage 等世界著名出版公司建立了良好的合作关系，从他们现有的数百种教材中甄选出 Andrew S. Tanenbaum, Bjarne Stroustrup, Brain W. Kernighan, Dennis Ritchie, Jim Gray, Alfred V. Aho, John E. Hopcroft, Jeffrey D. Ullman, Abraham Silberschatz, William Stallings, Donald E. Knuth, John L. Hennessy, Larry L. Peterson 等大师名家的一批经典作品，以“计算机科学丛书”为总称出版，供读者学习、研究及珍藏。大理石纹理的封面，也正体现了这套丛书的品位和格调。

“计算机科学丛书”的出版工作得到了国内外学者的鼎力襄助，国内的专家不仅提供了中肯的选题指导，还不辞劳苦地担任了翻译和审校的工作；而原书的作者也相当关注其作品在中国的传播，有的还专程为其书的中译本作序。迄今，“计算机科学丛书”已经出版了近两百个品种，这些书籍在读者中树立了良好的口碑，并被许多高校采用为正式教材和参考书籍。其影印版“经典原版书库”作为姊妹篇也被越来越多实施双语教学的学校所采用。

权威的作者、经典的教材、一流的译者、严格的审校、精细的编辑，这些因素使我们的图书有了质量的保证。随着计算机科学与技术专业学科建设的不断完善和教材改革的逐渐深化，教育界对国外计算机教材的需求和应用都将步入一个新的阶段，我们的目标是尽善尽美，而反馈的意见正是我们达到这一终极目标的重要帮助。华章公司欢迎老师和读者对我们的工作提出建议或给予指正，我们的联系方式如下：

华章网站：[www.hzbook.com](http://www.hzbook.com)

电子邮件：[hzsj@hzbook.com](mailto:hzsj@hzbook.com)

联系电话：(010) 88379604

联系地址：北京市西城区百万庄南街1号

邮政编码：100037



华章教育

华章科技图书出版中心

# Preface

Advances in data generation and collection are producing data sets of massive size in commerce and a variety of scientific disciplines. Data warehouses store details of the sales and operations of businesses, Earth-orbiting satellites beam high-resolution images and sensor data back to Earth, and genomics experiments generate sequence, structural, and functional data for an increasing number of organisms. The ease with which data can now be gathered and stored has created a new attitude toward data analysis: Gather whatever data you can whenever and wherever possible. It has become an article of faith that the gathered data will have value, either for the purpose that initially motivated its collection or for purposes not yet envisioned.

The field of data mining grew out of the limitations of current data analysis techniques in handling the challenges posed by these new types of data sets. Data mining does not replace other areas of data analysis, but rather takes them as the foundation for much of its work. While some areas of data mining, such as association analysis, are unique to the field, other areas, such as clustering, classification, and anomaly detection, build upon a long history of work on these topics in other fields. Indeed, the willingness of data mining researchers to draw upon existing techniques has contributed to the strength and breadth of the field, as well as to its rapid growth.

Another strength of the field has been its emphasis on collaboration with researchers in other areas. The challenges of analyzing new types of data cannot be met by simply applying data analysis techniques in isolation from those who understand the data and the domain in which it resides. Often, skill in building multidisciplinary teams has been as responsible for the success of data mining projects as the creation of new and innovative algorithms. Just as, historically, many developments in statistics were driven by the needs of agriculture, industry, medicine, and business, many of the developments in data mining are being driven by the needs of those same fields.

This book began as a set of notes and lecture slides for a data mining course that has been offered at the University of Minnesota since Spring 1998 to upper-division undergraduate and graduate students. Presentation slides

and notes developed in these offerings grew with time and served as a basis for the book. A survey of clustering techniques in data mining, originally written in preparation for research in the area, served as a starting point for one of the chapters in the book. Over time, the clustering chapter was joined by chapters on data, classification, association analysis, and anomaly detection. The book in its current form has been class tested at the home institutions of the authors—the University of Minnesota and Michigan State University—as well as several other universities.

A number of data mining books appeared in the meantime, but were not completely satisfactory for our students—primarily graduate and undergraduate students in computer science, but including students from industry and a wide variety of other disciplines. Their mathematical and computer backgrounds varied considerably, but they shared a common goal: to learn about data mining as directly as possible in order to quickly apply it to problems in their own domains. Thus, texts with extensive mathematical or statistical prerequisites were unappealing to many of them, as were texts that required a substantial database background. The book that evolved in response to these students' needs focuses as directly as possible on the key concepts of data mining by illustrating them with examples, simple descriptions of key algorithms, and exercises.

**Overview** Specifically, this book provides a comprehensive introduction to data mining and is designed to be accessible and useful to students, instructors, researchers, and professionals. Areas covered include data preprocessing, visualization, predictive modeling, association analysis, clustering, and anomaly detection. The goal is to present fundamental concepts and algorithms for each topic, thus providing the reader with the necessary background for the application of data mining to real problems. In addition, this book also provides a starting point for those readers who are interested in pursuing research in data mining or related fields.

The book covers five main topics: data, classification, association analysis, clustering, and anomaly detection. Except for anomaly detection, each of these areas is covered in a pair of chapters. For classification, association analysis, and clustering, the introductory chapter covers basic concepts, representative algorithms, and evaluation techniques, while the more advanced chapter discusses advanced concepts and algorithms. The objective is to provide the reader with a sound understanding of the foundations of data mining, while still covering many important advanced topics. Because of this approach, the book is useful both as a learning tool and as a reference.

To help the readers better understand the concepts that have been presented, we provide an extensive set of examples, figures, and exercises. Bibliographic notes are included at the end of each chapter for readers who are interested in more advanced topics, historically important papers, and recent trends. The book also contains a comprehensive subject and author index.

**To the Instructor** As a textbook, this book is suitable for a wide range of students at the advanced undergraduate or graduate level. Since students come to this subject with diverse backgrounds that may not include extensive knowledge of statistics or databases, our book requires minimal prerequisites—no database knowledge is needed and we assume only a modest background in statistics or mathematics. To this end, the book was designed to be as self-contained as possible. Necessary material from statistics, linear algebra, and machine learning is either integrated into the body of the text, or for some advanced topics, covered in the appendices.

Since the chapters covering major data mining topics are self-contained, the order in which topics can be covered is quite flexible. The core material is covered in Chapters 2, 4, 6, 8, and 10. Although the introductory data chapter (2) should be covered first, the basic classification, association analysis, and clustering chapters (4, 6, and 8, respectively) can be covered in any order. Because of the relationship of anomaly detection (10) to classification (4) and clustering (8), these chapters should precede Chapter 10. Various topics can be selected from the advanced classification, association analysis, and clustering chapters (5, 7, and 9, respectively) to fit the schedule and interests of the instructor and students. We also advise that the lectures be augmented by projects or practical exercises in data mining. Although they are time consuming, such hands-on assignments greatly enhance the value of the course.

**Support Materials** The supplements for the book are available at Addison-Wesley's Website [www.aw.com/cssupport](http://www.aw.com/cssupport). Support materials available to all readers of this book include

- PowerPoint lecture slides
- Suggestions for student projects
- Data mining resources such as data mining algorithms and data sets
- On-line tutorials that give step-by-step examples for selected data mining techniques described in the book using actual data sets and data analysis software

## viii Preface

Additional support materials, including solutions to exercises, are available only to instructors adopting this textbook for classroom use. Please contact your school's Addison-Wesley representative for information on obtaining access to this material. Comments and suggestions, as well as reports of errors, can be sent to the authors through [dmbook@cs.unm.edu](mailto:dmbook@cs.unm.edu).

**Acknowledgments** Many people contributed to this book. We begin by acknowledging our families to whom this book is dedicated. Without their patience and support, this project would have been impossible.

We would like to thank the current and former students of our data mining groups at the University of Minnesota and Michigan State for their contributions. Eui-Hong (Sam) Han and Mahesh Joshi helped with the initial data mining classes. Some of the exercises and presentation slides that they created can be found in the book and its accompanying slides. Students in our data mining groups who provided comments on drafts of the book or who contributed in other ways include Shyam Boriah, Haibin Cheng, Varun Chandola, Eric Eilertson, Levent Ertöz, Jing Gao, Rohit Gupta, Sridhar Iyer, Jung-Eun Lee, Benjamin Mayer, Aysel Ozgur, Uygur Oztekin, Gaurav Pandey, Kashif Riaz, Jerry Scripps, Gyorgy Simon, Hui Xiong, Jieping Ye, and Pusheng Zhang. We would also like to thank the students of our data mining classes at the University of Minnesota and Michigan State University who worked with early drafts of the book and provided invaluable feedback. We specifically note the helpful suggestions of Bernardo Craemer, Arifin Ruslim, Jamshid Vayghan, and Yu Wei.

Joydeep Ghosh (University of Texas) and Sanjay Ranka (University of Florida) class tested early versions of the book. We also received many useful suggestions directly from the following UT students: Pankaj Adhikari, Rajiv Bhatia, Frederic Bosche, Arindam Chakraborty, Meghana Deodhar, Chris Everson, David Gardner, Saad Godil, Todd Hay, Clint Jones, Ajay Joshi, Joonsoo Lee, Yue Luo, Anuj Nanavati, Tyler Olsen, Sunyoung Park, Aashish Phansalkar, Geoff Prewett, Michael Ryoo, Daryl Shannon, and Mei Yang.

Ronald Kostoff (ONR) read an early version of the clustering chapter and offered numerous suggestions. George Karypis provided invaluable L<sup>A</sup>T<sub>E</sub>X assistance in creating an author index. Irene Moulitsas also provided assistance with L<sup>A</sup>T<sub>E</sub>X and reviewed some of the appendices. Musetta Steinbach was very helpful in finding errors in the figures.

We would like to acknowledge our colleagues at the University of Minnesota and Michigan State who have helped create a positive environment for data mining research. They include Dan Boley, Joyce Chai, Anil Jain, Ravi

Janardan, Rong Jin, George Karypis, Haesun Park, William F. Punch, Shashi Shekhar, and Jaideep Srivastava. The collaborators on our many data mining projects, who also have our gratitude, include Ramesh Agrawal, Steve Cannon, Piet C. de Groen, Fran Hill, Yongdae Kim, Steve Klooster, Kerry Long, Nihar Mahapatra, Chris Potter, Jonathan Shapiro, Kevin Silverstein, Nevin Young, and Zhi-Li Zhang.

The departments of Computer Science and Engineering at the University of Minnesota and Michigan State University provided computing resources and a supportive environment for this project. ARDA, ARL, ARO, DOE, NASA, and NSF provided research support for Pang-Ning Tan, Michael Steinbach, and Vipin Kumar. In particular, Kamal Abdali, Dick Brackney, Jagdish Chandra, Joe Coughlan, Michael Coyle, Stephen Davis, Frederica Darema, Richard Hirsch, Chandrika Kamath, Raju Namburu, N. Radhakrishnan, James Sidoran, Bhavani Thuraisingham, Walt Tiernin, Maria Zemankova, and Xiaodong Zhang have been supportive of our research in data mining and high-performance computing.

It was a pleasure working with the helpful staff at Pearson Education. In particular, we would like to thank Michelle Brown, Matt Goldstein, Katherine Harutunian, Marilyn Lloyd, Kathy Smith, and Joyce Wells. We would also like to thank George Nichols, who helped with the art work and Paul Anagnostopoulos, who provided L<sup>A</sup>T<sub>E</sub>X support. We are grateful to the following Pearson reviewers: Chien-Chung Chan (University of Akron), Zhengxin Chen (University of Nebraska at Omaha), Chris Clifton (Purdue University), Joydeep Ghosh (University of Texas, Austin), Nazli Goharian (Illinois Institute of Technology), J. Michael Hardin (University of Alabama), James Hearne (Western Washington University), Hillol Kargupta (University of Maryland, Baltimore County and Agnik, LLC), Eamonn Keogh (University of California-Riverside), Bing Liu (University of Illinois at Chicago), Mariofanna Milanova (University of Arkansas at Little Rock), Srinivasan Parthasarathy (Ohio State University), Zbigniew W. Ras (University of North Carolina at Charlotte), Xintao Wu (University of North Carolina at Charlotte), and Mohammed J. Zaki (Rensselaer Polytechnic Institute).

# Contents

<b>Preface</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 What Is Data Mining? . . . . .	2
1.2 Motivating Challenges . . . . .	4
1.3 The Origins of Data Mining . . . . .	6
1.4 Data Mining Tasks . . . . .	7
1.5 Scope and Organization of the Book . . . . .	11
1.6 Bibliographic Notes . . . . .	13
1.7 Exercises . . . . .	16
<b>2 Data</b>	<b>19</b>
2.1 Types of Data . . . . .	22
2.1.1 Attributes and Measurement . . . . .	23
2.1.2 Types of Data Sets . . . . .	29
2.2 Data Quality . . . . .	36
2.2.1 Measurement and Data Collection Issues . . . . .	37
2.2.2 Issues Related to Applications . . . . .	43
2.3 Data Preprocessing . . . . .	44
2.3.1 Aggregation . . . . .	45
2.3.2 Sampling . . . . .	47
2.3.3 Dimensionality Reduction . . . . .	50
2.3.4 Feature Subset Selection . . . . .	52
2.3.5 Feature Creation . . . . .	55
2.3.6 Discretization and Binarization . . . . .	57
2.3.7 Variable Transformation . . . . .	63
2.4 Measures of Similarity and Dissimilarity . . . . .	65
2.4.1 Basics . . . . .	66
2.4.2 Similarity and Dissimilarity between Simple Attributes . . . . .	67
2.4.3 Dissimilarities between Data Objects . . . . .	69
2.4.4 Similarities between Data Objects . . . . .	72

2.4.5	Examples of Proximity Measures . . . . .	73
2.4.6	Issues in Proximity Calculation . . . . .	80
2.4.7	Selecting the Right Proximity Measure . . . . .	83
2.5	Bibliographic Notes . . . . .	84
2.6	Exercises . . . . .	88
<b>3</b>	<b>Exploring Data . . . . .</b>	<b>97</b>
3.1	The Iris Data Set . . . . .	98
3.2	Summary Statistics . . . . .	98
3.2.1	Frequencies and the Mode . . . . .	99
3.2.2	Percentiles . . . . .	100
3.2.3	Measures of Location: Mean and Median . . . . .	101
3.2.4	Measures of Spread: Range and Variance . . . . .	102
3.2.5	Multivariate Summary Statistics . . . . .	104
3.2.6	Other Ways to Summarize the Data . . . . .	105
3.3	Visualization . . . . .	105
3.3.1	Motivations for Visualization . . . . .	105
3.3.2	General Concepts . . . . .	106
3.3.3	Techniques . . . . .	110
3.3.4	Visualizing Higher-Dimensional Data . . . . .	124
3.3.5	Do's and Don'ts . . . . .	130
3.4	OLAP and Multidimensional Data Analysis . . . . .	131
3.4.1	Representing Iris Data as a Multidimensional Array . . . . .	131
3.4.2	Multidimensional Data: The General Case . . . . .	133
3.4.3	Analyzing Multidimensional Data . . . . .	135
3.4.4	Final Comments on Multidimensional Data Analysis . . . . .	139
3.5	Bibliographic Notes . . . . .	139
3.6	Exercises . . . . .	141
<b>4</b>	<b>Classification:</b>	
	<b>Basic Concepts, Decision Trees, and Model Evaluation . . . . .</b>	<b>145</b>
4.1	Preliminaries . . . . .	146
4.2	General Approach to Solving a Classification Problem . . . . .	148
4.3	Decision Tree Induction . . . . .	150
4.3.1	How a Decision Tree Works . . . . .	150
4.3.2	How to Build a Decision Tree . . . . .	151
4.3.3	Methods for Expressing Attribute Test Conditions . . . . .	155
4.3.4	Measures for Selecting the Best Split . . . . .	158
4.3.5	Algorithm for Decision Tree Induction . . . . .	164
4.3.6	An Example: Web Robot Detection . . . . .	166

## **xii Contents**

4.3.7	Characteristics of Decision Tree Induction . . . . .	168
4.4	Model Overfitting . . . . .	172
4.4.1	Overfitting Due to Presence of Noise . . . . .	175
4.4.2	Overfitting Due to Lack of Representative Samples . . .	177
4.4.3	Overfitting and the Multiple Comparison Procedure . .	178
4.4.4	Estimation of Generalization Errors . . . . .	179
4.4.5	Handling Overfitting in Decision Tree Induction . . . .	184
4.5	Evaluating the Performance of a Classifier . . . . .	186
4.5.1	Holdout Method . . . . .	186
4.5.2	Random Subsampling . . . . .	187
4.5.3	Cross-Validation . . . . .	187
4.5.4	Bootstrap . . . . .	188
4.6	Methods for Comparing Classifiers . . . . .	188
4.6.1	Estimating a Confidence Interval for Accuracy . . . . .	189
4.6.2	Comparing the Performance of Two Models . . . . .	191
4.6.3	Comparing the Performance of Two Classifiers . . . . .	192
4.7	Bibliographic Notes . . . . .	193
4.8	Exercises . . . . .	198
<b>5</b>	<b>Classification: Alternative Techniques</b>	<b>207</b>
5.1	Rule-Based Classifier . . . . .	207
5.1.1	How a Rule-Based Classifier Works . . . . .	209
5.1.2	Rule-Ordering Schemes . . . . .	211
5.1.3	How to Build a Rule-Based Classifier . . . . .	212
5.1.4	Direct Methods for Rule Extraction . . . . .	213
5.1.5	Indirect Methods for Rule Extraction . . . . .	221
5.1.6	Characteristics of Rule-Based Classifiers . . . . .	223
5.2	Nearest-Neighbor classifiers . . . . .	223
5.2.1	Algorithm . . . . .	225
5.2.2	Characteristics of Nearest-Neighbor Classifiers . . . . .	226
5.3	Bayesian Classifiers . . . . .	227
5.3.1	Bayes Theorem . . . . .	228
5.3.2	Using the Bayes Theorem for Classification . . . . .	229
5.3.3	Naïve Bayes Classifier . . . . .	231
5.3.4	Bayes Error Rate . . . . .	238
5.3.5	Bayesian Belief Networks . . . . .	240
5.4	Artificial Neural Network (ANN) . . . . .	246
5.4.1	Perceptron . . . . .	247
5.4.2	Multilayer Artificial Neural Network . . . . .	251
5.4.3	Characteristics of ANN . . . . .	255

5.5	Support Vector Machine (SVM) . . . . .	256
5.5.1	Maximum Margin Hyperplanes . . . . .	256
5.5.2	Linear SVM: Separable Case . . . . .	259
5.5.3	Linear SVM: Nonseparable Case . . . . .	266
5.5.4	Nonlinear SVM . . . . .	270
5.5.5	Characteristics of SVM . . . . .	276
5.6	Ensemble Methods . . . . .	276
5.6.1	Rationale for Ensemble Method . . . . .	277
5.6.2	Methods for Constructing an Ensemble Classifier . . . . .	278
5.6.3	Bias-Variance Decomposition . . . . .	281
5.6.4	Bagging . . . . .	283
5.6.5	Boosting . . . . .	285
5.6.6	Random Forests . . . . .	290
5.6.7	Empirical Comparison among Ensemble Methods . . . . .	294
5.7	Class Imbalance Problem . . . . .	294
5.7.1	Alternative Metrics . . . . .	295
5.7.2	The Receiver Operating Characteristic Curve . . . . .	298
5.7.3	Cost-Sensitive Learning . . . . .	302
5.7.4	Sampling-Based Approaches . . . . .	305
5.8	Multiclass Problem . . . . .	306
5.9	Bibliographic Notes . . . . .	309
5.10	Exercises . . . . .	315
<b>6</b>	<b>Association Analysis: Basic Concepts and Algorithms</b> . . . . .	<b>327</b>
6.1	Problem Definition . . . . .	328
6.2	Frequent Itemset Generation . . . . .	332
6.2.1	The <i>Apriori</i> Principle . . . . .	333
6.2.2	Frequent Itemset Generation in the <i>Apriori</i> Algorithm . . . . .	335
6.2.3	Candidate Generation and Pruning . . . . .	338
6.2.4	Support Counting . . . . .	342
6.2.5	Computational Complexity . . . . .	345
6.3	Rule Generation . . . . .	349
6.3.1	Confidence-Based Pruning . . . . .	350
6.3.2	Rule Generation in <i>Apriori</i> Algorithm . . . . .	350
6.3.3	An Example: Congressional Voting Records . . . . .	352
6.4	Compact Representation of Frequent Itemsets . . . . .	353
6.4.1	Maximal Frequent Itemsets . . . . .	354
6.4.2	Closed Frequent Itemsets . . . . .	355
6.5	Alternative Methods for Generating Frequent Itemsets . . . . .	359
6.6	FP-Growth Algorithm . . . . .	363

6.6.1	FP-Tree Representation . . . . .	363
6.6.2	Frequent Itemset Generation in FP-Growth Algorithm . . . . .	366
6.7	Evaluation of Association Patterns . . . . .	370
6.7.1	Objective Measures of Interestingness . . . . .	371
6.7.2	Measures beyond Pairs of Binary Variables . . . . .	382
6.7.3	Simpson's Paradox . . . . .	384
6.8	Effect of Skewed Support Distribution . . . . .	386
6.9	Bibliographic Notes . . . . .	390
6.10	Exercises . . . . .	404
<b>7</b>	<b>Association Analysis: Advanced Concepts . . . . .</b>	<b>415</b>
7.1	Handling Categorical Attributes . . . . .	415
7.2	Handling Continuous Attributes . . . . .	418
7.2.1	Discretization-Based Methods . . . . .	418
7.2.2	Statistics-Based Methods . . . . .	422
7.2.3	Non-discretization Methods . . . . .	424
7.3	Handling a Concept Hierarchy . . . . .	426
7.4	Sequential Patterns . . . . .	429
7.4.1	Problem Formulation . . . . .	429
7.4.2	Sequential Pattern Discovery . . . . .	431
7.4.3	Timing Constraints . . . . .	436
7.4.4	Alternative Counting Schemes . . . . .	439
7.5	Subgraph Patterns . . . . .	442
7.5.1	Graphs and Subgraphs . . . . .	443
7.5.2	Frequent Subgraph Mining . . . . .	444
7.5.3	<i>Apriori</i> -like Method . . . . .	447
7.5.4	Candidate Generation . . . . .	448
7.5.5	Candidate Pruning . . . . .	453
7.5.6	Support Counting . . . . .	457
7.6	Infrequent Patterns . . . . .	457
7.6.1	Negative Patterns . . . . .	458
7.6.2	Negatively Correlated Patterns . . . . .	458
7.6.3	Comparisons among Infrequent Patterns, Negative Pat- terns, and Negatively Correlated Patterns . . . . .	460
7.6.4	Techniques for Mining Interesting Infrequent Patterns . . . . .	461
7.6.5	Techniques Based on Mining Negative Patterns . . . . .	463
7.6.6	Techniques Based on Support Expectation . . . . .	465
7.7	Bibliographic Notes . . . . .	469
7.8	Exercises . . . . .	473

<b>8</b>	<b>Cluster Analysis: Basic Concepts and Algorithms</b>	<b>487</b>
8.1	Overview	490
8.1.1	What Is Cluster Analysis?	490
8.1.2	Different Types of Clusterings	491
8.1.3	Different Types of Clusters	493
8.2	K-means	496
8.2.1	The Basic K-means Algorithm	497
8.2.2	K-means: Additional Issues	506
8.2.3	Bisecting K-means	508
8.2.4	K-means and Different Types of Clusters	510
8.2.5	Strengths and Weaknesses	510
8.2.6	K-means as an Optimization Problem	513
8.3	Agglomerative Hierarchical Clustering	515
8.3.1	Basic Agglomerative Hierarchical Clustering Algorithm	516
8.3.2	Specific Techniques	518
8.3.3	The Lance-Williams Formula for Cluster Proximity	524
8.3.4	Key Issues in Hierarchical Clustering	524
8.3.5	Strengths and Weaknesses	526
8.4	DBSCAN	526
8.4.1	Traditional Density: Center-Based Approach	527
8.4.2	The DBSCAN Algorithm	528
8.4.3	Strengths and Weaknesses	530
8.5	Cluster Evaluation	532
8.5.1	Overview	533
8.5.2	Unsupervised Cluster Evaluation Using Cohesion and Separation	536
8.5.3	Unsupervised Cluster Evaluation Using the Proximity Matrix	542
8.5.4	Unsupervised Evaluation of Hierarchical Clustering	544
8.5.5	Determining the Correct Number of Clusters	546
8.5.6	Clustering Tendency	547
8.5.7	Supervised Measures of Cluster Validity	548
8.5.8	Assessing the Significance of Cluster Validity Measures	553
8.6	Bibliographic Notes	555
8.7	Exercises	559
<b>9</b>	<b>Cluster Analysis: Additional Issues and Algorithms</b>	<b>569</b>
9.1	Characteristics of Data, Clusters, and Clustering Algorithms	570
9.1.1	Example: Comparing K-means and DBSCAN	570
9.1.2	Data Characteristics	571

## xvi Contents

9.1.3	Cluster Characteristics . . . . .	573
9.1.4	General Characteristics of Clustering Algorithms . . . . .	575
9.2	Prototype-Based Clustering . . . . .	577
9.2.1	Fuzzy Clustering . . . . .	577
9.2.2	Clustering Using Mixture Models . . . . .	583
9.2.3	Self-Organizing Maps (SOM) . . . . .	594
9.3	Density-Based Clustering . . . . .	600
9.3.1	Grid-Based Clustering . . . . .	601
9.3.2	Subspace Clustering . . . . .	604
9.3.3	DENCLUE: A Kernel-Based Scheme for Density-Based Clustering . . . . .	608
9.4	Graph-Based Clustering . . . . .	612
9.4.1	Sparsification . . . . .	613
9.4.2	Minimum Spanning Tree (MST) Clustering . . . . .	614
9.4.3	OPOSSUM: Optimal Partitioning of Sparse Similarities Using METIS . . . . .	616
9.4.4	Chameleon: Hierarchical Clustering with Dynamic Modeling . . . . .	616
9.4.5	Shared Nearest Neighbor Similarity . . . . .	622
9.4.6	The Jarvis-Patrick Clustering Algorithm . . . . .	625
9.4.7	SNN Density . . . . .	627
9.4.8	SNN Density-Based Clustering . . . . .	629
9.5	Scalable Clustering Algorithms . . . . .	630
9.5.1	Scalability: General Issues and Approaches . . . . .	630
9.5.2	BIRCH . . . . .	633
9.5.3	CURE . . . . .	635
9.6	Which Clustering Algorithm? . . . . .	639
9.7	Bibliographic Notes . . . . .	643
9.8	Exercises . . . . .	647
10	Anomaly Detection . . . . .	651
10.1	Preliminaries . . . . .	653
10.1.1	Causes of Anomalies . . . . .	653
10.1.2	Approaches to Anomaly Detection . . . . .	654
10.1.3	The Use of Class Labels . . . . .	655
10.1.4	Issues . . . . .	656
10.2	Statistical Approaches . . . . .	658
10.2.1	Detecting Outliers in a Univariate Normal Distribution . . . . .	659
10.2.2	Outliers in a Multivariate Normal Distribution . . . . .	661
10.2.3	A Mixture Model Approach for Anomaly Detection . . . . .	662