



INTRODUCTION TO

DATA MINING



PANG-NING TAN
MICHAEL STEINBACH
VIPIN KUMAR

TP274
T161

INTRODUCTION TO **DATA MINING**

PANG-NING TAN

Michigan State University

MICHAEL STEINBACH

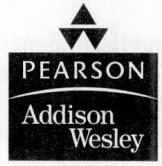
University of Minnesota

VIPIN KUMAR

University of Minnesota
and Army High Performance
Computing Research Center



E200601015



Boston San Francisco New York
London Toronto Sydney Tokyo Singapore Madrid
Mexico City Munich Paris Cape Town Hong Kong Montreal

Acquisitions Editor Matt Goldstein
Project Editor Katherine Harutunian
Production Supervisor Marilyn Lloyd
Production Services Paul C. Anagnostopoulos of Windfall Software
Marketing Manager Michelle Brown
Copyeditor Kathy Smith
Proofreader Jennifer McClain
Technical Illustration George Nichols
Cover Design Supervisor Joyce Cosentino Wells
Cover Design Night & Day Design
Cover Image © 2005 Rob Casey/Brand X Pictures
Prepress and Manufacturing Caroline Fell
Printer Hamilton Printing

Access the latest information about Addison-Wesley titles from our World Wide Web site:
<http://www.aw-bc.com/computing>

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and Addison-Wesley was aware of a trademark claim, the designations have been printed in initial caps or all caps.

The programs and applications presented in this book have been included for their instructional value. They have been tested with care, but are not guaranteed for any particular purpose. The publisher does not offer any warranties or representations, nor does it accept any liabilities with respect to the programs or applications.

Library of Congress Cataloging-in-Publication Data

A catalog record for this book is available from the Library of Congress.

Copyright © 2006 by Pearson Education, Inc.

For information on obtaining permission for use of material in this work, please submit a written request to Pearson Education, Inc., Rights and Contract Department, 75 Arlington Street, Suite 300, Boston, MA 02116 or fax your request to (617) 848-7047.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or any other media embodiments now known or hereafter to become known, without the prior written permission of the publisher. Printed in the United States of America.

ISBN 0-321-32136-7

2 3 4 5 6 7 8 9 10-HAM-08 07 06 05

INTRODUCTION TO **DATA MINING**

To our families ...

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 exercises 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. 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

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 dmbok@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^AT_EX assistance in creating an author index. Irene Moulitsas also provided assistance with L^AT_EX 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^AT_EX 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

| | |
|--|------------|
| Preface | vii |
| 1 Introduction | 1 |
| 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 |
| 2 Data | 19 |
| 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 |
| 3 | Exploring Data | 97 |
| 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 |
| 4 | Classification: | |
| | Basic Concepts, Decision Trees, and Model Evaluation | 145 |
| 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 |

| | | |
|----------|---|------------|
| 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 |
| 5 | Classification: Alternative Techniques | 207 |
| 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 |
| 6 | Association Analysis: Basic Concepts and Algorithms | 327 |
| 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 |
| 7 | Association Analysis: Advanced Concepts | 415 |
| 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>A priori</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 |

| | | |
|----------|---|------------|
| 8 | Cluster Analysis: Basic Concepts and Algorithms | 487 |
| 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 |
| 9 | Cluster Analysis: Additional Issues and Algorithms | 569 |
| 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 |

| | | |
|-----------|---|------------|
| 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 |