



# Introduction to Business Data Mining

DAVID OLSON  
YONG SHI

# Introduction to Business Data Mining

**David Olson**

*University of Nebraska–Lincoln*

**Yong Shi**

*Graduate University of the  
Chinese Academy of Sciences  
University of Nebraska–Omaha*

**Mc  
Graw  
Hill** **McGraw-Hill  
Irwin**

Boston Burr Ridge, IL Dubuque, IA Madison, WI New York San Francisco St. Louis  
Bangkok Bogotá Caracas Kuala Lumpur Lisbon London Madrid Mexico City  
Milan Montreal New Delhi Santiago Seoul Singapore Sydney Taipei Toronto



INTRODUCTION TO BUSINESS DATA MINING

Published by McGraw-Hill/Irwin, a business unit of The McGraw-Hill Companies, Inc., 1221 Avenue of the Americas, New York, NY, 10020. Copyright © 2007 by The McGraw-Hill Companies, Inc. All rights reserved. No part of this publication may be reproduced or distributed in any form or by any means, or stored in a database or retrieval system, without the prior written consent of The McGraw-Hill Companies, Inc., including, but not limited to, in any network or other electronic storage or transmission, or broadcast for distance learning.

Some ancillaries, including electronic and print components, may not be available to customers outside the United States.

This book is printed on acid-free paper.

1 2 3 4 5 6 7 8 9 0 CCW/CCW 0 9 8 7 6 5

ISBN-13: 978-0-07-295971-0

ISBN-10: 0-07-295971-1

Editorial director: *Brent Gordon*

Executive editor: *Scott Isenberg*

Editorial coordinator: *Lee Stone*

Senior marketing manager: *Douglas Reiner*

Senior media producer: *Victor Chiu*

Project manager: *Harvey Yep*

Senior production supervisor: *Sesha Bolisetty*

Designer: *Cara David*

Lead media project manager: *Brian Nacik*

Cover design: *Chris Bowyer*

Typeface: *10/12 Palatino*

Compositor: *International Typesetting & Composition*

Printer: *Courier Westford*

**Library of Congress Cataloging-in-Publication Data**

Olson, David Louis.

Introduction to business data mining / David Olson, Yong Shi.

p. cm.

Includes index.

ISBN-13: 978-0-07-295971-0 (alk. paper)

ISBN-10: 0-07-295971-1 (alk. paper)

1. Data mining. 2. Business—Data processing. I. Shi, Yong, 1956— II. Title.

HF5548.2.O46 2007

006.3'12—dc22

2005044525

To my family.

*David Olson*

This book is dedicated to my parents,  
Li Guihua and Shi Yuanqing for their constant  
support and understanding of my  
academic career.

*Yong Shi*

# About the Authors

## **David L. Olson**

*University of Nebraska–Lincoln*

**David Olson** is the James & H. K. Stuart Professor in MIS and Othmer Professor at the University of Nebraska. He received his Ph.D. in Business from the University of Nebraska in 1981. He has published research in over eighty refereed journal articles, primarily on the topic of multiple objective decision making. He teaches in the management information systems, management science, and operations management areas. He has authored the books *Decision Aids for Selection Problems*, *Introduction to Information Systems Project Management*, and *Managerial Issues of Enterprise Resource Planning Systems*, and co-authored the books *Decision Support Models and Expert Systems*; *Introduction to Management Science*; *Introduction to Simulation and Risk Analysis*; *Business Statistics: Quality Information for Decision Analysis*; *Statistics, Decision Analysis, and Decision Modeling*; *Multiple Criteria Analysis in Strategic Siting Problems*; and *Introduction to Business Data Mining*. He has made over one hundred presentations at international and national conferences on research topics. He is a member of the Association for Information Systems, the Decision Sciences Institute, the Institute for Operations Research and Management Sciences, and the Multiple Criteria Decision-Making Society. He has coordinated the Decision Sciences Institute Dissertation Competition and the Innovative Education Competition, chaired the Doctoral Affairs Committee, and served as nationally elected vice president three times and as National Program Chair. He was with Texas A&M University from 1981 through 2001—the last two years as Lowry Mays Professor of Business in the Department of Information and Operations Management. He received a Research Fellow Award from the College of Business and Graduate School of Business at Texas A&M University and held the Business Analysis Faculty Excellence Fellowship for two years. He is a Fellow of the Decision Sciences Institute.

## **Yong Shi**

*Graduate University of the Chinese Academy of Sciences  
University of Nebraska–Omaha*

**Professor Yong Shi** currently is the director of Chinese Academy of Sciences Research Center on Data Technology & Knowledge Economy and Assistant President of the Graduate University of Chinese Academy of Sciences. He has been the Charles W. and Margre H. Durham Distinguished Professor of Information Technology, College of Information Science and Technology, University of Nebraska at Omaha since 1999. Dr. Shi's research interests cover data mining, information overload, multiple criteria decision making, and telecommunications management. He has published seven books, more than sixty papers in various journals, and numerous conferences/proceedings papers. He is the Editor-in-Chief of *International Journal of Information Technology and Decision Making (SCI)*, an Area Editor of *International Journal of Operations and*

*Quantitative Management*, and an Editorial Board Member for a number of academic journals. Dr. Shi has received many distinguished awards including Outstanding Young Scientist Award, National Natural Science Foundation of China, 2001; Member of Overseas Assessor for the Chinese Academy of Sciences, May 2000; and Speaker of Distinguished Visitors Program (DVP) for 1997–2000, IEEE Computer Society. He has consulted for a number of famous companies in data mining and knowledge management projects.

# Preface

The intent of this book is to serve advanced undergraduate and graduate classes presenting data mining. Data mining is a very useful topic, applying quantitative analysis to large-scale data made available through recently developed information technology. Each of us has taught such material, and we both have extensive experience in quantitative analysis in business. Yong Shi also has extensive real experience in commercial data mining analysis. We want to take this opportunity to acknowledge the graduate students at the University of Nebraska at Omaha, Gang Kou, Nian Yan, and Wei Zhuang, who helped us prepare the data mining reports by using computer software for this book.

## Book Concept

Our intent is to cover the fundamental concepts of data mining, to demonstrate the potential of gathering large sets of data and analyzing these data sets to gain useful business understanding. We have organized the material into four parts. Part I introduces concepts. Part II describes and demonstrates basic data mining algorithms. Part III focuses on business applications of data mining. Part IV presents developing areas, including web mining, text mining, and ethical aspects of data mining. Part I is overview material. Part II contains chapters on a number of different techniques often used in data mining. Not all of these chapters need to be covered, and their sequence can be varied according to instructor need. Part III covers applications, and while we feel that these chapters contain the most interesting and important material, instructors who wish to focus on techniques might not wish to cover these chapters. Conversely, instructors more interested in business applications can cover Part III before reviewing content as needed in Part II. This approach would work especially well if data mining software is available to do the modeling. Part IV contains material we feel is important now and is growing in importance. However, again, coverage and sequence is up to the instructor.

The book includes short vignettes of how specific concepts have been applied in real business. A series of representative data sets are generated to demonstrate specific methods and concepts. References to data mining software and sites such as [www.kdnuggets.com](http://www.kdnuggets.com) are provided.

Supplements accompanying this text include (1) an instructor's CD-ROM, containing a solutions guide, PowerPoint slides, and the data set; (2) a student's CD-ROM, containing PowerPoint slides and the data set; and (3) an online learning center.

## Part I: Introduction

Chapter 1 gives an overview of data mining and provides a description of the data mining process. An overview of useful business applications is provided. Chapter 2 presents the data mining process in more detail. It demonstrates this process with a typical set of data. Visualization of data through data mining software is addressed. Chapter 3 presents database support to data mining. Different software tools are described, from data warehouse products through data marts to online analytic

processing. Data quality is addressed. Again, different concepts are demonstrated through prototypical data.

## **Part II: Data Mining Methods as Tools**

Chapter 4 provides an overview of data mining techniques and functions. Chapter 5 describes and demonstrates clustering algorithms. Software product support available is reviewed. Chapter 6 reviews various forms of regression tools to identify the best fit over given data sets. Chapter 7 discusses neural networks, a popular application of artificial intelligence suitable for many data mining applications. Chapter 8 reviews decision tree algorithms. The basic algorithm is described, along with descriptions of tree structure, machine learning, and fuzzy set aspects of decision trees. Software products are reviewed, and See5 is demonstrated. Chapter 9 presents linear programming-based methods of fitting data. Real data mining applications are described and demonstrated.

## **Part III: Business Applications**

Chapter 10 reviews the major applications of data mining in business, focusing on the value of these analyses to business decision making. This includes the important topics of customer relationship management. The concept of lift is described. The development of market segmentation by Fingerhut Inc. is reviewed. Chapter 11 describes market-basket analysis, a more qualitative data mining technique. This methodology is described through an example reported in the practitioner literature, and the fundamental data mining concepts of actionability, affinity positioning, and cross-selling are described.

## **Part IV: Developing Issues**

Chapter 12 presents text and web mining. Chapter 13 discusses ethical issues related to data mining.

*David L. Olson, University of Nebraska–Lincoln*

*Yong Shi, University of Nebraska–Omaha*



# Brief Table

## **PART ONE**

### **Introduction 1**

- 1** Initial Description of Data Mining in Business 3
- 2** Data Mining Processes and Knowledge Discovery 19
- 3** Database Support to Data Mining 34

## **PART TWO**

### **Data Mining Methods as Tools 51**

- 4** Overview of Data Mining Techniques 53
- 5** Cluster Analysis 73
- 6** Regression Algorithms in Data Mining 99
- 7** Neural Networks in Data Mining 122

- 8** Decision Tree Algorithms 135

- 9** Linear Programming–Based Methods 164

## **PART THREE**

### **Business Applications 187**

- 10** Business Data Mining Applications 189
- 11** Market-Basket Analysis 211

## **PART FOUR**

### **Developing Issues 223**

- 12** Text and Web Mining 225
- 13** Ethical Aspects of Data Mining 250

## **GLOSSARY 261**

## **INDEX 265**

# Contents

## **PART ONE**

### **INTRODUCTION 1**

#### **Chapter 1**

#### **Initial Description of Data Mining in Business 3**

- Introduction 4
- What Is Needed to Do Data Mining 5
- Data Mining 5
- Focused Marketing 7
- Business Data Mining 8
  - Retailing* 8
  - Banking* 9
  - Credit Card Management* 9
  - Insurance* 10
  - Telecommunications* 11
  - Telemarketing* 12
  - Human Resource Management* 13
- Data Mining Tools 13
- Summary 14

#### **Chapter 2**

#### **Data Mining Processes and Knowledge Discovery 19**

- CRISP-DM 20
  - Business Understanding* 21
  - Data Understanding* 21
  - Data Preparation* 22
  - Modeling* 24
  - Evaluation* 26
  - Deployment* 27
- Knowledge Discovery Process 27
- Summary 31

#### **Chapter 3**

#### **Database Support to Data Mining 34**

- Data Warehousing 35
- Data Marts 36
- Online Analytic Processing 37
- Data Warehouse Implementation 38
- Metadata 40
- System Demonstrations 41

- Data Warehouse* 41
- Data Mart* 42
- OLAP* 42
- Data Quality 44
- Software Products 45
- Real Examples 46
  - Wal-Mart's Data Warehouse System* 46
  - Summers Rubber Company Data Storage Design* 46
- Summary 48

## **PART TWO**

### **DATA MINING METHODS**

### **AS TOOLS 51**

#### **Chapter 4**

#### **Overview of Data Mining Techniques 53**

- Data Mining Models 54
- Data Mining Perspectives 55
- Data Mining Functions 56
- Demonstration Data Sets 57
  - Loan Application Data* 58
  - Job Application Data* 59
  - Insurance Fraud Data* 60
  - Expenditure Data* 62
- Appendix: Enterprise Miner Demonstration on Expenditure Data Set 63
  - Data Partitioning* 63
  - Regression Modeling* 64
  - Decision Tree Modeling* 68
  - Neural Network Modeling* 69
- Summary 70

#### **Chapter 5**

#### **Cluster Analysis 73**

- Cluster Analysis 74
- Description of Cluster Analysis 74
  - A Clustering Algorithm* 75
  - Insurance Fraud Data* 75
  - Weighted Distance Cluster Model* 79
- Varying the Number of Clusters 79
  - The Three-Cluster Model* 83
- Applications of Cluster Analysis 83
  - Monitoring Credit Card Accounts* 84
  - Data Mining of Insurance Claims* 84

Clustering Methods Used in Software 85  
 Application of Methods to Larger Data Sets 86  
     *Loan Application Data* 86  
     *Insurance Fraud Data* 87  
     *Expenditure Data* 89  
 Software Products 91  
 Appendix: Clementine 91  
     *Web Plot* 92  
 Summary 96

## Chapter 6

### Regression Algorithms in Data Mining 99

Regression Models 100  
     *Classical Tests of the Regression Model* 103  
     *Multiple Regression* 104  
 Logistic Regression 107  
 Linear Discriminant Analysis 109  
     *Discriminant Function for Loan Data* 109  
     *Job Applicant Data* 110  
 Real Applications of Regression in Data Mining 113  
     *Stepwise Regression in Bankruptcy Prediction Models* 113  
 Application of Models to Larger Data Sets 114  
     *Insurance Fraud Data* 114  
     *Job Applicant Data* 117  
     *Loan Applicant Data* 118  
 Summary 119

## Chapter 7

### Neural Networks in Data Mining 122

Neural Networks 123  
     *An Example Neural Network Application* 124  
 Neural Networks in Data Mining 126  
 Business Applications of Neural Networks 126  
     *Neural Network Models for Bankruptcy Prediction* 126  
     *Data Mining to Target Customers* 127  
 Application of Neural Networks to Larger Data Sets 128  
     *Insurance Fraud Data* 128  
     *Job Applicant Data* 129  
     *Loan Applicant Data* 130  
 Neural Network Products 130  
 Summary 131

## Chapter 8

### Decision Tree Algorithms 135

Decision Tree Operation 136  
     *Rule Interestingness* 137

Machine Learning 138  
 Decision Tree Applications 143  
     *Inventory Prediction* 143  
     *The Mining of Clinical Databases* 144  
     *Software Development Quality* 144  
     *Evaluation* 146  
 Application of Methods to Larger Data Sets 147  
     *Loan Application Data* 147  
     *Insurance Fraud Data* 151  
     *Job Application Data* 152  
 Decision Tree Software Products 153  
 Appendix: Demonstration of See5 Decision Tree Analysis 154  
     *Data Cleaning* 154  
     *Data Mining Process* 155  
 Summary 160

## Chapter 9

### Linear Programming–Based Methods 164

Linear Discriminant Analysis 165  
 Multiple Criteria Linear Programming Classification 168  
 Fuzzy Linear Programming Classification 172  
 Credit Card Portfolio Management: A Real-Life Application 176  
 Linear Programming–Based Software Support 182  
 Appendix: Data Mining Linear Programming Formulations 182  
     *Multiple-Class Separation* 182  
     *Multiple Criteria Linear Programming Separation* 183  
     *Fuzzy Linear Programming Separation* 184  
 Summary 184

## PART THREE

### BUSINESS APPLICATIONS 187

## Chapter 10

### Business Data Mining Applications 189

Applications 190  
     *Mailstream Optimization at Fingerhut* 191  
     *Customer Relationship Management (CRM)* 194  
     *Credit Scoring* 196  
     *Investment Risk Analysis* 200  
     *Data Mining Applications in Insurance* 202  
 Comparisons of Data Mining Methods 206  
 Summary 208

**Chapter 11****Market-Basket Analysis 211**

- Definitions 212
- Demonstration 214
- Market-Basket Limitations 216
- Market-Basket Analysis Software 217
- Appendix: Market-Basket Procedure 217
- Summary 219

**PART FOUR****DEVELOPING ISSUES 223****Chapter 12****Text and Web Mining 225**

- Text Mining 226
  - Identification of Key Terms* 226
  - Text Mining Demonstration* 227
  - Link Analysis* 228
  - Text Mining Applied to Legal Databases* 230
  - Text Mining Products* 231
- Web Mining 232
  - Web Mining Taxonomy* 233
  - Web User Behavior* 234
  - Web Mining Examples* 235

*Web Mining Systems* 236

*Web Usage Mining Demonstration* 2236

**Appendix: Semantic Text Analysis 237**

*Semantic Text Analysis* 238

*Quantitative Models* 240

*Unstructured Text Analysis* 242

*Taxonomy Classification* 242

*OLAP Charts* 246

*Summary* 246

Summary 247

**Chapter 13****Ethical Aspects of Data Mining 250**

The Hazards of Data Access 251

Web Data Mining Issues 253

The Problem 254

Web Ethics 254

*Privacy Issues* 255

*Discrimination* 256

Methods of Control 256

Summary 257

**Glossary 261****Index 265**

PART ONE

# Introduction



# Initial Description of Data Mining in Business

This chapter:

- Introduces data mining concepts
- Presents typical business data mining applications
- Explains the meaning of key concepts
- Gives a brief overview of data mining tools
- Outlines the remaining chapters of the book

Our culture has developed the ability to generate masses of data. Computer systems expand much faster than the human ability to absorb. Furthermore, Internet connections make it possible to share data in real time on a global basis.

Recent political events emphasize the existence of data to predict. Some blame the system when terrorist strikes are not prevented, because if you dig deep enough, you can always find some data or a memo that pointed to the coming occurrence of these events. However, you would also find a great deal more data predicting things that didn't happen. Obviously, there's a clear need for many organizations to be able to process data faster and more reliably. Data mining involves the use of analysis to detect patterns and allow predictions. Though it's not a perfect science, the intent of data mining is to gain small advantages, because perfect predictions are impossible. These small advantages can be extremely profitable to business. For instance, retail sales organizations have developed sophisticated customer segmentation models to save them from sending sales materials to consumers who likely won't purchase their products, focusing instead on those segments with a higher probability of sales. Banks and other organizations have developed sophisticated customer relationship management programs (supported by data mining) that can predict the value of specific types of customers to that organization, and predict repayment of loans as well. Insurance companies have long applied statistical analysis, which has been extended by data mining tools to aid in the prediction of fraudulent claims. These are only three of many important data mining applications to business.

This book seeks to describe some business applications of data mining. It also will describe the general process of data mining, those database tools needed to support data mining, and the techniques available for data mining.

## Introduction

Data mining refers to the analysis of the large quantities of data that are stored in computers. For example, grocery stores have large amounts of data generated by our purchases. Bar coding has made checking out very convenient for us, and provides retail establishments with masses of data. Grocery stores and other retail stores are able to quickly process our purchases, and use computers to accurately determine product prices. These same computers can help the stores with their inventory management, by instantaneously determining the quantity of items of each product on hand. They are also able to apply computer technology to contact their vendors so they don't run out of the things that we want to purchase. Computers allow the store's accounting system to more accurately measure costs, and determine the profit that store stockholders are concerned about. All of this information is available based upon the bar code information attached to each product. Along with many other sources of information, data gathered through bar coding can be used for data mining analysis.

Data mining is not limited to business. Both major parties in the 2004 U.S. election utilized data mining of potential voters.<sup>1</sup> Data mining has been heavily used in the medical field, to include patient diagnosis records to help identify best practices.<sup>2</sup> The Mayo Clinic worked with IBM to develop an online computer system to identify how that last 100 Mayo patients with the same gender, age, and medical history responded to particular treatments.<sup>3</sup>

Business use of data mining is also impressive. Toyota used the data mining of its **data warehouse** to determine more efficient transportation routes, reducing the time to deliver cars to customers by an average of 19 days. Data warehouses (to be discussed in Chapter 3) are enormous database systems capable of systematically storing all transactional data generated by a business organization, such as WalMart. Toyota also was able to identify sales trends faster, and identify the best locations for new dealerships. Benefits were estimated to be \$30 million per year in North America.<sup>4</sup>

Data mining is widely used by banking firms in soliciting credit card customers,<sup>5</sup> by insurance and telecommunication companies in detecting fraud,<sup>6</sup> by manufacturing firms in quality control,<sup>7</sup> and many other applications. Data mining is being applied to improve food product safety,<sup>8</sup> criminal detection,<sup>9</sup> and tourism.<sup>10</sup> Fingerhut has become very successful in **micromarketing**, targeting small groups of highly responsive customers. Media companies such as R. R. Donnelly & Sons provide consumer and life-style data, as well as customized individual publications to firms that use data mining for catalog marketing.

Data mining involves statistical and/or artificial intelligence analysis, usually applied to large-scale data sets. Traditional statistical analysis involves an approach that is usually directed, in that a specific set of expected outcomes exists. This approach is referred to as supervised. However, there is more to data mining than the technical tools used. Data mining involves a spirit of **knowledge discovery** (learning new and useful things), which is referred to as unsupervised. Much of this can be accomplished through automatic means, as we will see in decision tree analysis, for example. But data mining is not limited to automated analysis. Knowledge discovery by humans can be enhanced by graphical tools and the identification of unexpected patterns through a combination of human and computer interaction.



Data mining can be used by businesses in many ways. Three examples are

- **Customer profiling** Identifying those subsets of customers that are most profitable to the business
- **Targeting** Determining the characteristics of profitable customers who have been captured by competitors
- **Market-basket analysis** Determining product purchases by consumers, which can be used for product positioning and for cross-selling.

These are not the only applications of data mining, but they are three of the most important to businesses. Customer profiling is a key part of customer relationship management (CRM), which will be elaborated upon in Chapter 10. Targeting is a key concept in managing **churn**, or customer turnover, also discussed in Chapter 10. Market-basket analysis is an interesting use of data mining that we discuss in Chapter 11.

## What Is Needed to Do Data Mining?

Data mining requires the identification of a problem, along with collection of data that can lead to a better understanding of the market, and computer models to provide statistical or other means of analysis. There are two general types of data mining studies. **Hypothesis testing** involves expressing a theory about the relationship between actions and outcomes. In a simple form, it can be hypothesized that advertising will yield greater profit. This relationship has long been studied by retailing firms in the context of their specific operations. Data mining is applied to identifying relationships based on large quantities of data, which could include testing the response rates to various types of advertising on the sales and profitability of specific product lines. The second form of data mining study is knowledge discovery. In this form of analysis, a preconceived notion may not be present, but rather relationships can be identified by looking at the data. This may be supported by visualization tools that display data, or through fundamental statistical analysis, such as correlation analysis.

A variety of analytic computer models have been used in data mining. Chapters 5 through 9 of this book will discuss various types of these models. Also required is access to data. Quite often, systems including data warehouses and data marts are used to manage large quantities of data (see Chapter 3). Other data mining analyses are done with smaller sets of data, such as can be organized in online analytic processing systems.

## Data Mining

Data mining has been called exploratory data analysis, among other things. Masses of data generated from cash registers, from scanning, and from topic-specific databases throughout the company are explored, analyzed, reduced, and reused. Searches are performed across different models proposed for predicting sales, marketing response, and profit. Classical statistical approaches are fundamental to data mining. Automated AI methods are also used. However, systematic exploration through classical statistical methods is still the basis of data mining. Some of the tools developed by the field of statistical analysis are harnessed through automatic control (with some key human guidance) in dealing with data.