# DATA MINING FOR BUSINESS ANALYTICS

## CONCEPTS, TECHNIQUES, AND APPLICATIONS IN R

Galit Shmueli • Peter C. Bruce
Inbal Yahav • Nitin R. Patel
Kenneth C. Lichtendahl Jr.

with website

WILEY

# DATA MINING
# FOR BUSINESS ANALYTICS

# DATA MINING
# FOR BUSINESS ANALYTICS

Concepts, Techniques, and Applications in R

GALIT SHMUELI

PETER C. BRUCE

INBAL YAHAV

NITIN R. PATEL

KENNETH C. LICHTENDAHL, JR.

WILEY

*The beginning of wisdom is this:*
*Get wisdom, and whatever else you get, get insight.*

רֵאשִׁית חָכְמָה, קְנֵה חָכְמָה;    וּבְכָל-קִנְיָנְךָ, קְנֵה בִינָה.

– Proverbs 4:7

# Contents

## PART II    DATA EXPLORATION AND DIMENSION REDUCTION

## CHAPTER 3    Data Visualization                                                 55

## CHAPTER 4    Dimension Reduction                                                 91

## Part III  PERFORMANCE EVALUATION

## Chapter 5  Evaluating Predictive Performance                              117

# Part IV PREDICTION AND CLASSIFICATION METHODS

## Chapter 6 Multiple Linear Regression     153

## Chapter 7 $k$-Nearest Neighbors ($k$NN)     173

## Chapter 8 The Naive Bayes Classifier     187

## PART V  MINING RELATIONSHIPS AMONG RECORDS

## CHAPTER 14 Association Rules and Collaborative Filtering    327

## CHAPTER 15 Cluster Analysis    355

## PART VI  FORECASTING TIME SERIES

### CHAPTER 16  Handling Time Series                                        385

### CHAPTER 17  Regression-Based Forecasting                                 399

# PART VIII CASES

## CHAPTER 21 Cases                                                                          497

# Foreword by Gareth James

The field of statistics has existed in one form or another for 200 years, and by the second half of the 20th century had evolved into a well-respected and essential academic discipline. However, its prominence expanded rapidly in the 1990s with the explosion of new, and enormous, data sources. For the first part of this century, much of this attention was focused on biological applications, in particular, genetics data generated as a result of the sequencing of the human genome. However, the last decade has seen a dramatic increase in the availability of data in the business disciplines, and a corresponding interest in business-related statistical applications.

The impact has been profound. Ten years ago, when I was able to attract a full class of MBA students to my new statistical learning elective, my colleagues were astonished because our department struggled to fill most electives. Today, we offer a Masters in Business Analytics, which is the largest specialized masters program in the school and has application volume rivaling those of our MBA programs. Our department's faculty size and course offerings have increased dramatically, yet the MBA students are still complaining that the classes are all full. Google's chief economist, Hal Varian, was indeed correct in 2009 when he stated that "the sexy job in the next 10 years will be statisticians."

This demand is driven by a simple, but undeniable, fact. Business analytics solutions have produced significant and measurable improvements in business performance, on multiple dimensions and in numerous settings, and as a result, there is a tremendous demand for individuals with the requisite skill set. However, training students in these skills is challenging given that, in addition to the obvious required knowledge of statistical methods, they need to understand business-related issues, possess strong communication skills, and be comfortable dealing with multiple computational packages. Most statistics texts concentrate on abstract training in classical methods, without much emphasis on practical, let alone business, applications.

This book has by far the most comprehensive review of business analytics methods that I have ever seen, covering everything from classical approaches such as linear and logistic regression, through to modern methods like neural

networks, bagging and boosting, and even much more business specific proce-
dures such as social network analysis and text mining. If not the bible, it is at
the least a definitive manual on the subject. However, just as important as the
list of topics, is the way that they are all presented in an applied fashion using
business applications. Indeed the last chapter is entirely dedicated to 10 separate
cases where business analytics approaches can be applied.

In this latest edition, the authors have added an important new dimension
in the form of the R software package. Easily the most widely used and influ-
ential open source statistical software, R has become the go-to tool for such
purposes. With literally hundreds of freely available add-on packages, R can
be used for almost any business analytics related problem. The book provides
detailed descriptions and code involving applications of R in numerous business
settings, ensuring that the reader will actually be able to apply their knowledge
to real-life problems.

We recently introduced a business analytics course into our required MBA
core curriculum and I intend to make heavy use of this book in developing the
syllabus. I'm confident that it will be an indispensable tool for any such course.

GARETH JAMES

*Marshall School of Business, University of Southern California, 2017*

# Foreword by Ravi Bapna

Data is the new gold—and mining this gold to create business value in today's context of a highly networked and digital society requires a skillset that we haven't traditionally delivered in business or statistics or engineering programs on their own. For those businesses and organizations that feel overwhelmed by today's Big Data, the phrase *you ain't seen nothing yet* comes to mind. Yesterday's three major sources of Big Data—the 20+ years of investment in enterprise systems (ERP, CRM, SCM, …), the 3 billion plus people on the online social grid, and the close to 5 billion people carrying increasingly sophisticated mobile devices—are going to be dwarfed by tomorrow's smarter physical ecosystems fueled by the Internet of Things (IoT) movement.

The idea that we can use sensors to connect physical objects such as homes, automobiles, roads, even garbage bins and streetlights, to digitally optimized systems of governance goes hand in glove with bigger data and the need for deeper analytical capabilities. We are not far away from a smart refrigerator sensing that you are short on, say, eggs, populating your grocery store's mobile app's shopping list, and arranging a Task Rabbit to do a grocery run for you. Or the refrigerator negotiating a deal with an Uber driver to deliver an evening meal to you. Nor are we far away from sensors embedded in roads and vehicles that can compute traffic congestion, track roadway wear and tear, record vehicle use and factor these into dynamic usage-based pricing, insurance rates, and even taxation. This brave new world is going to be fueled by analytics and the ability to harness data for competitive advantage.

Business Analytics is an emerging discipline that is going to help us ride this new wave. This new Business Analytics discipline requires individuals who are grounded in the fundamentals of business such that they know the right questions to ask, who have the ability to harness, store, and optimally process vast datasets from a variety of structured and unstructured sources, and who can then use an array of techniques from machine learning and statistics to uncover new insights for decision-making. Such individuals are a rare commodity today, but their creation has been the focus of this book for a decade now. This book's forte is that it relies on explaining the core set of concepts required for today's business analytics professionals using real-world data-rich cases in a hands-on manner,

without sacrificing academic rigor. It provides a modern day foundation for Business Analytics, the notion of linking the x's to the y's of interest in a predictive sense. I say this with the confidence of someone who was probably the first adopter of the zeroth edition of this book (Spring 2006 at the Indian School of Business).

I can't say enough about the long-awaited R edition. R is my go-to platform for analytics these days. It's also used by a wide variety of instructors in our MS-Business Analytics program. The open-innovation paradigm used by R is one key part of the analytics perfect storm, the other components being the advances in computing and the business appetite for data-driven decision-making.

I look forward to using the book in multiple fora, in executive education, in MBA classrooms, in MS-Business Analytics programs, and in Data Science bootcamps. I trust you will too!

RAVI BAPNA

*Carlson School of Management, University of Minnesota, 2017*

# Preface to the R Edition

This textbook first appeared in early 2007 and has been used by numerous students and practitioners and in many courses, ranging from dedicated data mining classes to more general business analytics courses (including our own experience teaching this material both online and in person for more than 10 years). The first edition, based on the Excel add-in XLMiner, was followed by two more XLMiner editions, a JMP edition, and now this R edition, with its companion website, www.dataminingbook.com.

This new R edition, which relies on the free and open-source R software, presents output from R, as well as the code used to produce that output, including specification of a variety of packages and functions. Unlike computer-science or statistics-oriented textbooks, the focus in this book is on data mining concepts, and how to implement the associated algorithms in R. We assume a basic facility with R.

For this R edition, two new co-authors stepped on board—Inbal Yahav and Casey Lichtendahl—bringing both expertise teaching business analytics courses using R and data mining consulting experience in business and government. Such practical experience is important, since the open-source nature of R software makes available a plethora of approaches, packages, and functions available for data mining. Given the main goal of this book—to introduce data mining concepts using R software for illustration—our challenge was to choose an R code cocktail that supports highlighting the important concepts. In addition to providing R code and output, this edition also incorporates updates and new material based on feedback from instructors teaching MBA, undergraduate, diploma, and executive courses, and from their students as well.

One update, compared to the first two editions of the book, is the title: we now use *Business Analytics* in place of *Business Intelligence*. This reflects the change in terminology since the second edition: Business Intelligence today refers mainly to reporting and data visualization ("what is happening now"), while Business Analytics has taken over the "advanced analytics," which include predictive analytics and data mining. In this new edition, we therefore use the updated terms.

This R edition includes the material that was recently added in the third edition of the original (XLMiner-based) book:

- Social network analysis

- Text mining

- Ensembles

- Uplift modeling

- Collaborative filtering

Since the appearance of the (XLMiner-based) second edition, the landscape of the courses using the textbook has greatly expanded: whereas initially, the book was used mainly in semester-long elective MBA-level courses, it is now used in a variety of courses in Business Analytics degrees and certificate programs, ranging from undergraduate programs, to post-graduate and executive education programs. Courses in such programs also vary in their duration and coverage. In many cases, this textbook is used across multiple courses. The book is designed to continue supporting the general "Predictive Analytics" or "Data Mining" course as well as supporting a set of courses in dedicated business analytics programs.

A general "Business Analytics," "Predictive Analytics," or "Data Mining" course, common in MBA and undergraduate programs as a one-semester elective, would cover Parts I–III, and choose a subset of methods from Parts IV and V. Instructors can choose to use cases as team assignments, class discussions, or projects. For a two-semester course, Part VI might be considered, and we recommend introducing the new Part VII (Data Analytics).

For a set of courses in a dedicated business analytics program, here are a few courses that have been using our book:

**Predictive Analytics: Supervised Learning** In a dedicated Business Analytics program, the topic of Predictive Analytics is typically instructed across a set of courses. The first course would cover Parts I–IV and instructors typically choose a subset of methods from Part IV according to the course length. We recommend including the new Chapter 13 in such a course, as well as the new "Part VII: Data Analytics."

**Predictive Analytics: Unsupervised Learning** This course introduces data exploration and visualization, dimension reduction, mining relationships, and clustering (Parts III and V). If this course follows the Predictive Analytics: Supervised Learning course, then it is useful to examine examples and approaches that integrate unsupervised and supervised learning, such as the new part on "Data Analytics."

**Forecasting Analytics** A dedicated course on time series forecasting would rely on Part VI.

**Advanced Analytics** A course that integrates the learnings from Predictive Analytics (supervised and unsupervised learning). Such a course can focus on Part VII: Data Analytics, where social network analytics and text mining are introduced. Some instructors choose to use the Cases (Chapter 21) in such a course.

In all courses, we strongly recommend including a project component, where data are either collected by students according to their interest or provided by the instructor (e.g., from the many data mining competition datasets available). From our experience and other instructors' experience, such projects enhance the learning and provide students with an excellent opportunity to understand the strengths of data mining and the challenges that arise in the process.

# Acknowledgments

We thank the many people who assisted us in improving the first three editions of the initial XLMiner version of this book and the JMP edition, as well as those who helped with comments on early drafts of this R edition. Anthony Babinec, who has been using earlier editions of this book for years in his data mining courses at Statistics.com, provided us with detailed and expert corrections. Dan Toy and John Elder IV greeted our project with early enthusiasm and provided detailed and useful comments on initial drafts. Ravi Bapna, who used an early draft in a data mining course at the Indian School of Business, has provided invaluable comments and helpful suggestions since the book's start.

Many of the instructors, teaching assistants, and students using earlier editions of the book have contributed invaluable feedback both directly and indirectly, through fruitful discussions, learning journeys, and interesting data mining projects that have helped shape and improve the book. These include MBA students from the University of Maryland, MIT, the Indian School of Business, National Tsing Hua University, and Statistics.com. Instructors from many universities and teaching programs, too numerous to list, have supported and helped improve the book since its inception.

Several professors have been especially helpful with this R edition: Hayri Tongarlak, Prashant Joshi (UKA Tarsadia University), Jay Annadatha, Roger Bohn, and Sridhar Vaithianathan provided detailed comments and R code files for the companion website; Scott Nestler has been a helpful friend of this book project from the beginning.

Kuber Deokar, instructional operations supervisor at Statistics.com, has been unstinting in his assistance, support, and detailed attention. We also thank Shweta Jadhav and Dhanashree Vishwasrao, assistant teachers. Valerie Troiano has shepherded many instructors and students through the Statistics.com courses that have helped nurture the development of these books.

Colleagues and family members have been providing ongoing feedback and assistance with this book project. Boaz Shmueli and Raquelle Azran gave detailed editorial comments and suggestions on the first two editions; Bruce McCullough and Adam Hughes did the same for the first edition. Noa Shmueli provided careful proofs of the third edition. Ran Shenberger offered design tips.