

UseR!

Chris Chapman
Elea McDonnell Feit

R for Marketing Research and Analytics

 Springer

Use R!

Series Editors:

Robert Gentleman Kurt Hornik Giovanni Parmigiani

More information about this series at <http://www.springer.com/series/6991>

Use R!

- Kolaczyk / Csárdi*: Statistical Analysis of Network Data with R (2014)
- Nolan / Temple Lang*: XML and Web Technologies for Data Sciences with R (2014)
- Willekens*: Multistate Analysis of Life Histories with R (2014)
- Cortez*: Modern Optimization with R (2014)
- Eddelbuettel*: Seamless R and C++ Integration with Rcpp (2013)
- Bivand / Pebesma / Gómez-Rubio*: Applied Spatial Data Analysis with R
(2nd ed. 2013)
- van den Boogaart / Tolosana-Delgado*: Analyzing Compositional Data with R
(2013)
- Nagarajan / Scutari / Lèbre*: Bayesian Networks in R (2013)

Chris Chapman • Elea McDonnell Feit

R for Marketing Research and Analytics

 Springer

Chris Chapman
Google, Inc.
Seattle, WA, USA
cnchapman+r@gmail.com

Elea McDonnell Feit
LeBow College of Business
Drexel University
Philadelphia, PA, USA
efeit@drexel.edu

ISSN 2197-5736

Use R!

ISBN 978-3-319-14435-1

DOI 10.1007/978-3-319-14436-8

ISSN 2197-5744 (electronic)

ISBN 978-3-319-14436-8 (eBook)

Library of Congress Control Number: 2014960277

Springer Cham Heidelberg New York Dordrecht London

© Springer International Publishing Switzerland 2015

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

Springer International Publishing AG Switzerland is part of Springer Science+Business Media
(www.springer.com)

Praise for R for Marketing Research and Analytics

R for Marketing Research and Analytics is the perfect book for those interested in driving success for their business and for students looking to get an introduction to R. While many books take a purely academic approach, Chapman (Google) and Feit (formerly of GM and the Modellers) know exactly what is needed for practical marketing problem solving. I am an expert R user, yet had never thought about a textbook that provides the soup-to-nuts way that Chapman and Feit do: show how to load a data set, explore it using visualization techniques, analyze it using statistical models, and then demonstrate the business implications. It is a book that I wish I had written.

Eric Bradlow, K.P. Chao Professor, Chairperson, Wharton Marketing Department and Co-Director, Wharton Customer Analytics Initiative

R for Marketing Research and Analytics provides an excellent introduction to the R statistical package for marketing researchers. This is a must-have book for anyone who seriously pursues analytics in the field of marketing. R is the software gold standard in the research industry, and this book provides an introduction to R and shows how to run the analysis. Topics range from graphics and exploratory methods to confirmatory methods including structural equation modeling, all illustrated with data. A great contribution to the field!

Greg Allenby, Helen C. Kurtz Chair in Marketing, Professor of Marketing, Professor of Statistics, Ohio State University

Chris Chapman's and Elea Feit's engaging and authoritative book nicely fills a gap in the literature. At last we have an accessible book that presents core marketing research methods using the tools and vernacular of modern data science. The book will enable marketing researchers to up their game by adopting the R statistical computing environment. And data scientists with an interest in marketing problems now have a reference that speaks to them in their language.

James Guszczka, Chief Data Scientist, Deloitte Consulting – US

Finally a highly accessible guide for getting started with R. Feit and Chapman have applied years of lessons learned to developing this easy-to-use guide, designed to quickly build a strong foundation for applying R to sound analysis. The authors succeed in demystifying R by employing a likeable and practical writing style, along with sensible organization and comfortable pacing of the material. In addition to covering all the most important analysis techniques, the authors are generous throughout in providing tips for optimizing R's efficiency and identifying common pitfalls. With this guide, anyone interested in R can begin using it confidently in a short period of time for analysis, visualization, and for more advanced analytics procedures. *R for Marketing Research and Analytics* is the perfect guide and reference text for the casual and advanced user alike.

Matt Valle, Executive Vice President, Global Key Account Management – GfK

Preface

We are here to help you learn R for marketing research and analytics.

R is a great choice for marketing analysts. It offers unsurpassed capabilities for fitting statistical models. It is extensible and is able to process data from many different systems, in a variety of forms, for both small and large data sets. The R ecosystem includes the widest available range of established and emerging statistical methods as well as visualization techniques. Yet the use of R in marketing lags other fields such as statistics, econometrics, psychology, and bioinformatics. With your help, we hope to change that!

This book is designed for two audiences: practicing marketing researchers and analysts who want to learn R, and students or researchers from other fields who want to review selected marketing topics in an R context.

What are the prerequisites? Simply that you are interested in R for marketing, are conceptually familiar with basic statistical models such as linear regression, and are willing to engage in hands-on learning. This book will be particularly helpful to analysts who have some degree of programming experience and wish to learn R. In Chap. 1 we describe additional reasons to use R (and a few reasons perhaps *not* to use R).

The *hands-on* part is important. We teach concepts gradually in a sequence across the first seven chapters and ask you to *type* our examples as you work; this book is *not* a cookbook-style reference. We spend some time (as little as possible) in Part I on the basics of the R language and then turn in Part II to applied, real-world marketing analytics problems. Part III presents a few advanced marketing topics. Every chapter shows off the power of R, and we hope each one will teach you something new and interesting.

Specific features of this book are as follows:

- It is organized around marketing research tasks. Instead of generic examples, we put methods into the context of marketing questions.

- We presume only basic statistics knowledge and use a minimum of mathematics. This book is designed to be approachable for practitioners and does not dwell on equations or mathematical details of statistical models (although we give references to those texts).
- This is a didactic book that explains statistical concepts and the R code. We want you to understand what we're doing and learn how to avoid common problems in both statistics and R. We intend the book to be *readable* and to fulfill a different need than references and cookbooks available elsewhere.
- The applied chapters demonstrate progressive model building. We do not present “the answer” but instead show how an analyst might realistically conduct analyses in successive steps where multiple models are compared for statistical strength and practical utility.
- The chapters include visualization as a part of core analyses. We don't regard visualization as a stand-alone topic; rather, we believe it is an integral part of data exploration and model building.
- You will learn more than just R. In addition to core models, we include topics such as structural models and transaction analysis that may be new and useful even for experienced analysts.
- The book reflects both traditional and Bayesian approaches. Core models are presented with traditional (frequentist) methods, while later sections introduce Bayesian methods for linear models and conjoint analysis.
- Most of the analyses use simulated data, which provides practice in the R language along with additional insight into the structure of marketing data. If you are inclined, you can change the data simulation and see how the statistical models are affected.
- Where appropriate, we call out more advanced material on programming or models so that you may either skip it or read it, as you find appropriate. These sections are indicated by * in their titles (such as *This is an advanced section**).

What do we *not* cover? For one, this book teaches *R* for marketing and does not teach marketing research in itself. We discuss many marketing topics but omit others that would simply repeat the analytic methods in R. As noted above, we approach statistical models from a conceptual point of view and skip the mathematics. A few specialized topics have been omitted due to complexity and space; these include customer lifetime value models and econometric time series models. Overall, we believe the analyses here represent a great sample of marketing research and analytics practice. If you learn to perform these, you'll be well equipped to apply R in many areas of marketing.

Why are we the right teachers? We've used R and its predecessor S for a combined 27 years since 1997 and it is our primary analytics platform. We perform marketing analyses of all kinds in R, ranging from simple data summaries to complex analyses involving thousands of lines of custom code and newly created models.

We've also taught R to many people. This book grew from courses the authors have presented at American Marketing Association (AMA) events including the Academy of Marketing Analytics at Emory University and several years of the Advanced Research Techniques Forum (ART Forum). We have also taught R at the Sawtooth Software Conference and to students and industry collaborators at the Wharton School. We thank those many students for their feedback and believe that their experiences will benefit you.

Acknowledgements

We want to give special thanks here to people who made this book possible. First are all the students from our tutorials and classes over the years. They provided valuable feedback, and we hope their experiences will benefit you.

In the marketing academic and practitioner community, we had valuable feedback from Ken Deal, Fred Feinberg, Shane Jensen, Jake Lee, Dave Lyon, and Bruce McCullough.

Chris's colleagues in the research community at Google provided extensive feedback on portions of the book. We thank Mario Callegaro, Marianna Dizik, Rohan Gifford, Tim Hesterberg, Shankar Kumar, Norman Lemke, Paul Litvak, Katrina Panovich, Marta Rey-Babarro, Kerry Rodden, Dan Russell, Angela Schörgendorfer, Steven Scott, Bob Silverstein, Gill Ward, John Webb, and Yori Zwols for their encouragement and comments.

The staff and editors at Springer helped us smooth the process, especially Hannah Bracken, Jon Gurstelle, and the Use R! series editors.

Much of this book was written in public and university libraries, and we thank them for their hospitality alongside their unsurpassed literary resources. Portions of the book were written during pleasant days at the New Orleans Public Library, New York Public Library, Christoph Keller Jr. Library at the General Theological Seminary in New York, University of California San Diego Geisel Library, University of Washington Suzzallo and Allen Libraries, Sunnyvale Public Library, and most particularly, where the first words, code, and outline were written, along with much more later, the Tokyo Metropolitan Central Library.

Our families supported us in weekends and nights of editing, and they endured more discussion of R than is fair for any layperson. Thank you, Cristi, Maddie, Jeff, and Zoe.

Most importantly, we thank *you*, the reader. We're glad you've decided to investigate R, and we hope to repay your effort. Let's start!

New York, NY and Seattle, WA
Philadelphia, PA
November 2014

Chris Chapman
Elea McDonnell Feit

Contents

Preface	vii
----------------------	-----

Part I Basics of R

1 Welcome to R	3
1.1 What Is R?	3
1.2 Why R?	4
1.3 Why Not R?	5
1.4 When R?	6
1.5 Using This Book	6
1.5.1 About the Text	6
1.5.2 About the Data	7
1.5.3 Online Material	8
1.5.4 When Things Go Wrong	9
1.6 Key Points	10
2 An Overview of the R Language	11
2.1 Getting Started	11
2.1.1 Initial Steps	11
2.1.2 Starting R	12
2.2 A Quick Tour of R's Capabilities	13
2.3 Basics of Working with R Commands	17
2.4 Basic Objects	18
2.4.1 Vectors	19
2.4.2 Help! A Brief Detour	21
2.4.3 More on Vectors and Indexing	24
2.4.4 aaRgh! A Digression for New Programmers	26
2.4.5 Missing and Interesting Values	26
2.4.6 Using R for Mathematical Computation	28
2.4.7 Lists	28

- 2.5 Data Frames 30
- 2.6 Loading and Saving Data 34
 - 2.6.1 Image Files 36
 - 2.6.2 CSV Files 36
- 2.7 Writing Your Own Functions* 38
 - 2.7.1 Language Structures* 40
 - 2.7.2 Anonymous Functions* 41
- 2.8 Clean Up! 42
- 2.9 Learning More* 43
- 2.10 Key Points 44

Part II Fundamentals of Data Analysis

- 3 Describing Data** 47
 - 3.1 Simulating Data 47
 - 3.1.1 Store Data: Setting the Structure 48
 - 3.1.2 Store Data: Simulating Data Points 50
 - 3.2 Functions to Summarize a Variable 52
 - 3.2.1 Discrete Variables 52
 - 3.2.2 Continuous Variables 54
 - 3.3 Summarizing Data Frames 56
 - 3.3.1 `summary()` 57
 - 3.3.2 `describe()` 58
 - 3.3.3 Recommended Approach to Inspecting Data 59
 - 3.3.4 `apply() *` 59
 - 3.4 Single Variable Visualization 61
 - 3.4.1 Histograms 61
 - 3.4.2 Boxplots 66
 - 3.4.3 QQ Plot to Check Normality* 68
 - 3.4.4 Cumulative Distribution* 69
 - 3.4.5 Language Brief: `by()` and `aggregate()` 70
 - 3.4.6 Maps 72
 - 3.5 Learning More* 74
 - 3.6 Key Points 75
- 4 Relationships Between Continuous Variables** 77
 - 4.1 Retailer Data 77
 - 4.1.1 Simulating Customer Data 78
 - 4.1.2 Simulating Online and In-Store Sales Data 79
 - 4.1.3 Simulating Satisfaction Survey Responses 80
 - 4.1.4 Simulating Non-Response Data 82
 - 4.2 Exploring Associations Between Variables with Scatterplots 83
 - 4.2.1 Creating a Basic Scatterplot with `plot()` 83
 - 4.2.2 Color-Coding Points on a Scatterplot 86

4.2.3	Adding a Legend to a Plot	88
4.2.4	Plotting on a Log Scale	89
4.3	Combining Plots in a Single Graphics Object	90
4.4	Scatterplot Matrices	92
4.4.1	<code>pairs()</code>	92
4.4.2	<code>scatterplotMatrix()</code>	93
4.5	Correlation Coefficients	95
4.5.1	Correlation Tests	97
4.5.2	Correlation Matrices	98
4.5.3	Transforming Variables before Computing Correlations	100
4.5.4	Typical Marketing Data Transformations	102
4.5.5	Box–Cox Transformations*	102
4.6	Exploring Associations in Survey Responses*	104
4.6.1	<code>jitter()*</code>	105
4.6.2	<code>polychoric()*</code>	106
4.7	Learning More*	107
4.8	Key Points	108
5	Comparing Groups: Tables and Visualizations	111
5.1	Simulating Consumer Segment Data	111
5.1.1	Segment Data Definition	112
5.1.2	Language Brief: <code>for()</code> Loops	114
5.1.3	Language Brief: <code>if()</code> Blocks	116
5.1.4	Final Segment Data Generation	118
5.2	Finding Descriptives by Group	120
5.2.1	Language Brief: Basic Formula Syntax	123
5.2.2	Descriptives for Two-Way Groups	124
5.2.3	Visualization by Group: Frequencies and Proportions	126
5.2.4	Visualization by Group: Continuous Data	129
5.3	Learning More*	132
5.4	Key Points	133
6	Comparing Groups: Statistical Tests	135
6.1	Data for Comparing Groups	135
6.2	Testing Group Frequencies: <code>chisq.test()</code>	136
6.3	Testing Observed Proportions: <code>binom.test()</code>	139
6.3.1	About Confidence Intervals	140
6.3.2	More About <code>binom.test()</code> and Binomial Distributions	141
6.4	Testing Group Means: <code>t.test()</code>	142
6.5	Testing Multiple Group Means: ANOVA	144
6.5.1	Model Comparison in ANOVA*	146
6.5.2	Visualizing Group Confidence Intervals	147
6.5.3	Variable Selection in ANOVA: Stepwise Modeling*	148
6.6	Bayesian ANOVA: Getting Started*	149
6.6.1	Why Bayes?	150

- 6.6.2 Basics of Bayesian ANOVA* 150
- 6.6.3 Inspecting the Posterior Draws* 152
- 6.6.4 Plotting the Bayesian Credible Intervals* 155
- 6.7 Learning More* 156
- 6.8 Key Points 157
- 7 Identifying Drivers of Outcomes: Linear Models** 159
- 7.1 Amusement Park Data 160
 - 7.1.1 Simulating the Amusement Park Data 160
- 7.2 Fitting Linear Models with `lm()` 162
 - 7.2.1 Preliminary Data Inspection 163
 - 7.2.2 Recap: Bivariate Association 165
 - 7.2.3 Linear Model with a Single Predictor 165
 - 7.2.4 `lm` Objects 166
 - 7.2.5 Checking Model Fit 169
- 7.3 Fitting Linear Models with Multiple Predictors 173
 - 7.3.1 Comparing Models 175
 - 7.3.2 Using a Model to Make Predictions 176
 - 7.3.3 Standardizing the Predictors 177
- 7.4 Using Factors as Predictors 179
- 7.5 Interaction Terms 182
 - 7.5.1 Language Brief: Advanced Formula Syntax* 183
- 7.6 Caution! Overfitting 185
- 7.7 Recommended Procedure for Linear Model Fitting 186
- 7.8 Bayesian Linear Models with `MCMCregress()`* 186
- 7.9 Learning More* 188
- 7.10 Key Points 190

Part III Advanced Marketing Applications

- 8 Reducing Data Complexity** 195
- 8.1 Consumer Brand Rating Data 195
 - 8.1.1 Rescaling the Data 197
 - 8.1.2 Aggregate Mean Ratings by Brand 198
- 8.2 Principal Component Analysis and Perceptual Maps 200
 - 8.2.1 PCA Example 200
 - 8.2.2 Visualizing PCA 203
 - 8.2.3 PCA for Brand Ratings 204
 - 8.2.4 Perceptual Map of the Brands 206
 - 8.2.5 Cautions with Perceptual Maps 208
- 8.3 Exploratory Factor Analysis 209
 - 8.3.1 Basic EFA Concepts 210
 - 8.3.2 Finding an EFA Solution 211

8.3.3	EFA Rotations	213
8.3.4	Using Factor Scores for Brands	216
8.4	Multidimensional Scaling	218
8.4.1	Non-metric MDS	219
8.5	Learning More*	221
8.5.1	Principal Component Analysis	221
8.5.2	Factor Analysis	221
8.5.3	Multidimensional Scaling	222
8.6	Key Points	222
8.6.1	Principal Component Analysis	222
8.6.2	Exploratory Factor Analysis	222
8.6.3	Multidimensional Scaling	223
9	Additional Linear Modeling Topics	225
9.1	Handling Highly Correlated Variables	226
9.1.1	An Initial Linear Model of Online Spend	226
9.1.2	Remediating Collinearity	229
9.2	Linear Models for Binary Outcomes: Logistic Regression	231
9.2.1	Basics of the Logistic Regression Model	231
9.2.2	Data for Logistic Regression of Season Passes	232
9.2.3	Sales Table Data	233
9.2.4	Language Brief: Classes and Attributes of Objects*	234
9.2.5	Finalizing the Data	236
9.2.6	Fitting a Logistic Regression Model	237
9.2.7	Reconsidering the Model	239
9.2.8	Additional Discussion	242
9.3	Hierarchical Linear Models	242
9.3.1	Some HLM Concepts	243
9.3.2	Ratings-Based Conjoint Analysis for the Amusement Park	244
9.3.3	Simulating Ratings-Based Conjoint Data	245
9.3.4	An Initial Linear Model	246
9.3.5	Hierarchical Linear Model with <code>lme4</code>	248
9.3.6	The Complete Hierarchical Linear Model	249
9.3.7	Summary of HLM with <code>lme4</code>	251
9.4	Bayesian Hierarchical Linear Models*	252
9.4.1	Initial Linear Model with <code>MCMCregress()</code> *	253
9.4.2	Hierarchical Linear Model with <code>MCMChregress()</code> *	253
9.4.3	Inspecting Distribution of Preference*	256
9.5	A Quick Comparison of Frequentist & Bayesian HLMs*	259
9.6	Learning More*	263
9.6.1	Collinearity	263
9.6.2	Logistic Regression	263
9.6.3	Hierarchical Models	263
9.6.4	Bayesian Hierarchical Models	263
9.7	Key Points	264
9.7.1	Collinearity	264

- 9.7.2 Logistic Regression 264
- 9.7.3 Hierarchical Linear Models 265
- 9.7.4 Bayesian Methods for Hierarchical Linear Models 266
- 10 Confirmatory Factor Analysis and Structural Equation Modeling ... 267**
 - 10.1 The Motivation for Structural Models 268
 - 10.1.1 Structural Models in This Chapter 269
 - 10.2 Scale Assessment: CFA 270
 - 10.2.1 Simulating PIES CFA Data 272
 - 10.2.2 Estimating the PIES CFA Model 277
 - 10.2.3 Assessing the PIES CFA Model 278
 - 10.3 General Models: Structural Equation Models 283
 - 10.3.1 The Repeat Purchase Model in R 284
 - 10.3.2 Assessing the Repeat Purchase Model 286
 - 10.4 The Partial Least Squares (PLS) Alternative 288
 - 10.4.1 PLS-SEM for Repeat Purchase 289
 - 10.4.2 Visualizing the Fitted PLS Model* 292
 - 10.4.3 Assessing the PLS-SEM Model 293
 - 10.4.4 PLS-SEM with the Larger Sample 295
 - 10.5 Learning More* 297
 - 10.6 Key Points 297
- 11 Segmentation: Clustering and Classification 299**
 - 11.1 Segmentation Philosophy 299
 - 11.1.1 The Difficulty of Segmentation 299
 - 11.1.2 Segmentation as Clustering and Classification 300
 - 11.2 Segmentation Data 302
 - 11.3 Clustering 302
 - 11.3.1 The Steps of Clustering 303
 - 11.3.2 Hierarchical Clustering: `hclust()` Basics 305
 - 11.3.3 Hierarchical Clustering Continued: Groups from `hclust()` 309
 - 11.3.4 Mean-Based Clustering: `kmeans()` 311
 - 11.3.5 Model-Based Clustering: `Mclust()` 314
 - 11.3.6 Comparing Models with `BIC()` 315
 - 11.3.7 Latent Class Analysis: `poLCA()` 317
 - 11.3.8 Comparing Cluster Solutions 320
 - 11.3.9 Recap of Clustering 322
 - 11.4 Classification 322
 - 11.4.1 Naive Bayes Classification: `naiveBayes()` 323
 - 11.4.2 Random Forest Classification: `randomForest()` 327
 - 11.4.3 Random Forest Variable Importance 330
 - 11.5 Prediction: Identifying Potential Customers* 333
 - 11.6 Learning More* 336
 - 11.7 Key Points 337

12 Association Rules for Market Basket Analysis	339
12.1 The Basics of Association Rules	340
12.1.1 Metrics	340
12.2 Retail Transaction Data: Market Baskets	341
12.2.1 Example Data: Groceries	342
12.2.2 Supermarket Data	344
12.3 Finding and Visualizing Association Rules	346
12.3.1 Finding and Plotting Subsets of Rules	348
12.3.2 Using Profit Margin Data with Transactions: An Initial Start	349
12.3.3 Language Brief: A Function for Margin Using an Object's <code>class</code> *	351
12.4 Rules in Non-Transactional Data: Exploring Segments Again	356
12.4.1 Language Brief: Slicing Continuous Data with <code>cut()</code>	356
12.4.2 Exploring Segment Associations	357
12.5 Learning More*	360
12.6 Key Points	360
13 Choice Modeling	363
13.1 Choice-Based Conjoint Analysis Surveys	364
13.2 Simulating Choice Data*	365
13.3 Fitting a Choice Model	370
13.3.1 Inspecting Choice Data	371
13.3.2 Fitting Choice Models with <code>mlogit()</code>	372
13.3.3 Reporting Choice Model Findings	375
13.3.4 Share Predictions for Identical Alternatives	380
13.3.5 Planning the Sample Size for a Conjoint Study	381
13.4 Adding Consumer Heterogeneity to Choice Models	383
13.4.1 Estimating Mixed Logit Models with <code>mlogit()</code>	383
13.4.2 Share Prediction for Heterogeneous Choice Models	386
13.5 Hierarchical Bayes Choice Models	388
13.5.1 Estimating Hierarchical Bayes Choice Models with <code>ChoiceModelR</code>	388
13.5.2 Share Prediction for Hierarchical Bayes Choice Models	395
13.6 Design of Choice-Based Conjoint Surveys*	397
13.7 Learning More*	398
13.8 Key Points	399
Conclusion	401
A Appendix: R Versions and Related Software	403
A.1 R Base	403
A.2 RStudio	404
A.3 Emacs Speaks Statistics	405
A.4 Eclipse + StatET	406
A.5 Revolution R	407

A.6	Other Options	408
A.6.1	Text Editors	408
A.6.2	R Commander	408
A.6.3	Rattle	409
A.6.4	Deducer	409
A.6.5	TIBCO Enterprise Runtime for R	409
B	Appendix: Scaling Up	411
B.1	Handling Data	411
B.1.1	Data Wrangling	411
B.1.2	Microsoft Excel: <code>gdata</code>	412
B.1.3	SAS, SPSS, and Other Statistics Packages: <code>foreign</code>	412
B.1.4	SQL: <code>RSQLite</code> , <code>sqldf</code> and <code>RODBC</code>	413
B.2	Handling Large Data Sets	415
B.3	Speeding Up Computation	416
B.3.1	Efficient Coding and Data Storage	416
B.3.2	Enhancing the R Engine	417
B.4	Time Series Analysis, Repeated Measures, and Longitudinal Analysis	418
B.5	Automated and Interactive Reporting	419
C	Appendix: Packages Used	423
C.1	Core and Frequentist Statistics	424
C.2	Graphics	424
C.3	Bayesian Methods	425
C.4	Advanced Statistics	426
C.5	Machine Learning	426
C.6	Data Handling	427
C.7	Other Packages	428
D	Appendix: Online Materials and Data Files	431
D.1	Data File Structure	431
D.2	Data File URL Cross-Reference	432
D.2.1	Update on Data Locations	432
	References	435
	Index	447

Part I

Basics of R

Welcome to R

1.1 What Is R?

As a marketing analyst, you have no doubt heard of R. You may have tried R and become frustrated and confused, after which you returned to other tools that are “good enough.” You may know that R uses a command line and dislike that. Or you may be convinced of R’s advantages for experts but worry that you don’t have time to learn or use it.

We are here to help! Our goal is to present *just the essentials*, in the *minimal necessary time*, with *hands-on learning* so you will come up to speed as quickly as possible to be productive in R. In addition, we’ll cover a few advanced topics that demonstrate the power of R and might teach advanced users some new skills.

A key thing to realize is that *R is a programming language*. It is *not* a “statistics program” like SPSS, SAS, JMP, or Minitab, and doesn’t wish to be one. The official R Project describes R as “a language and environment for statistical computing and graphics.” Notice that “language” comes first, and that “statistical” is coequal with “graphics.” R is a great programming language for doing statistics. The inventor of the underlying language, John Chambers received the 1998 Association for Computing Machinery (ACM) Software System Award for a system that “will forever alter the way people analyze, visualize, and manipulate data . . .”[6].

R was based on Chambers’s preceding S language (S as in “statistics”) developed in the 1970s and 1980s at Bell Laboratories, home of the UNIX operating system and the C programming language. S gained traction among analysts and academics in the 1990s as implemented in a commercial software package, S-PLUS. Robert Gentleman and Ross Ihaka wished to make the S approach more widely available and offered R as an open source project starting in 1997.

Since then, the popularity of R has grown geometrically. The real magic of R is that its users are able to contribute developments that enhance R with everything from additional core functions to highly specialized methods. And many do contribute! Today there are over 6,000 packages of add-on functionality available for R (see <http://cran.r-project.org/web/packages> for the latest count).

If you have experience in programming, you will appreciate some of R's key features right away. If you're new to programming, this chapter describes why R is special and Chap. 2 introduces the fundamentals of programming in R.

1.2 Why R?

There are many reasons to learn and use R. It is the platform of choice for the largest number of statisticians who create new analytics methods, so emerging techniques are often available first in R. R is rapidly becoming the default educational platform in university statistics programs and is spreading to other disciplines such as economics and psychology.

For analysts, R offers the largest and most diverse set of analytic tools and statistical methods. It allows you to write analyses that can be reused and that extend the R system itself. It runs on most operating systems and interfaces well with data systems such as online data and SQL databases. R offers beautiful and powerful plotting functions that are able to produce graphics vastly more tailored and informative than typical spreadsheet charts. Putting all of those together, R can vastly improve an analyst's overall productivity. Elea knows an enterprising analyst who used R to automate the process of downloading data and producing a formatted monthly report. The automation saved him almost 40 h of work each month . . . which he didn't tell his manager for a few months!

Then there is the community. Many R users are enthusiasts who love to help others and are rewarded in turn by the simple joy of solving problems and the fact that they often learn something new. R is a dynamic system created by its users, and there is always something new to learn. Knowledge of R is a valuable skill in demand for analytics jobs at a growing number of top companies.

R code is also inspectable; you may choose to trust it, yet you are also free to verify. All of its core code and most packages that people contribute are open source. You can examine the code to see exactly how analyses work and what is happening under the hood.

Finally, R is free. It is a labor of love and professional pride for the R Core Development Team, which includes eminent statisticians and computer scientists. As with all masterpieces, the quality of their devotion is evident in the final work.

1.3 Why Not R?

What's not to love? No doubt you've observed that not everyone in the world uses R. Being R-less is unimaginable to us, yet there are reasons why some analysts might not want to use it.

One reason not to use R is this: until you've mastered the basics of the language, many simple analyses are cumbersome to do in R. If you're new to R and want a table of means, cross-tabs, or a t-test, it may be frustrating to figure out how to get them. R is about power, flexibility, control, iterative analyses, and cutting-edge methods, not point-and-click deliverables.

Another reason is if you do not like programming. If you're new to programming, R is a great place to start. But if you've tried programming before and didn't enjoy it, R will be a challenge as well. Our job is to help you as much as we can, and we will try hard to teach R to you. However, not everyone enjoys programming. On the other hand, if you're an experienced coder, R will seem simple (perhaps deceptively so), and we will help you avoid a few pitfalls.

Some companies and their information technology or legal departments are skeptical of R because it is open source. It is common for managers to ask, "If it's free, how can it be good?" There are many responses to that, including pointing out the hundreds of books on R, its citation in peer-reviewed articles, and the list of eminent contributors (in R, run the `contributors()` command and web search some of them). Or you might try the engineer's adage: "It can be good, fast, or cheap: pick 2." R is good and cheap, but not fast, insofar as it requires time and effort to master.

As for R being free, you should realize that contributors to R actually do derive benefit; it just happens to be non-monetary. They are compensated through respect and reputation, through the power their own work gains, and by the contributions back to the ecosystem from other users. This is a rational economic model even when the monetary price is zero.

A final concern about R is the unpredictability of its ecosystem. With packages contributed by thousands of authors, there are priceless contributions along with others that are mediocre or flawed. The downside of having access to the latest developments is that many will not stand the test of time. It is up to you to determine whether a method meets your needs, and you cannot always rely on curation or authorities to determine it for you (although you will rapidly learn which authors and which experts' recommendations to trust). If you trust your judgment, this situation is no different than with any software. *Caveat emptor.*

We hope to convince you that for many purposes, the benefits of R outweigh the difficulties.

1.4 When R?

There are a few common use cases for R:

- You want access to methods that are newer or more powerful than available elsewhere. Many R users start for exactly that reason; they see a method in a journal article, conference paper, or presentation, and discover that the method is available only in R.
- You need to run an analysis many, many times. This is how Chris started his R journey; for his dissertation, he needed to bootstrap existing methods in order to compare their typical results to those of a new machine learning model. R is perfect for model iteration.
- You need to apply an analysis to multiple data sets. Because everything is scripted, R is great for analyses that are repeated across data sets. It even has tools available for automated reporting.
- You need to develop a new analytic technique or wish to have perfect control and insight into an existing method. For many statistical procedures, R is easier to code than other programming languages.
- Your manager, professor, or coworker is encouraging you to use R. We've influenced students and colleagues in this way and are happy to report that a large number of them are enthusiastic R users today.

By showing you the power of R, we hope to convince you that your current tools are *not* perfectly satisfactory. Even more deviously, we hope to rewrite your expectations about what *is* satisfactory.

1.5 Using This Book

This book is intended to be *didactic* and *hands-on*, meaning that we want to teach you about R and the models we use in plain English, and we expect you to engage with the code interactively in R. It is designed for you to type the commands as you read. (We also provide code files for download from the book's website; see Sect. 1.5.3 below.)

1.5.1 About the Text

R commands for you to run are presented in code blocks like this:

```
> citation()  
To cite R in publications use:
```

```
R Core Team (2014). R: A language and environment for statistical computing.
R Foundation for Statistical
Computing, Vienna, Austria. URL http://www.R-project.org/.
...
```

We describe these code blocks and interacting with R in Chap. 2. The code generally follows the Google style guide for R (available at <http://google-style-guide.googlecode.com/svn/trunk/Rguide.xml>) except when we thought a deviation might make the code or text clearer. (As you learn R, you will wish to make your code readable; the Google guide is very useful for code formatting.)

When we refer to R commands, add-on packages, or data in the text outside of code blocks, we set the names in monospace type like this: `citation()`. We include parentheses on function (command) names to indicate that they are functions, such as the `summary()` function (Sect. 2.4.1), as opposed to an object such as the `Groceries` data set (Sect. 12.2.1).

When we introduce or define significant new concepts, we set them in italic, such as *vectors*. Italic is also used simply for *emphasis*.

We teach the R language progressively throughout the book, and much of our coverage of the language is blended into chapters that cover marketing topics and statistical models. In those cases, we present crucial language topics in *Language Brief* sections (such as Sect. 3.4.5). To learn as much as possible about the R language, you'll need to read the Language Brief sections even if you only skim the surrounding material on statistical models.

Some sections cover deeper details or more advanced topics, and may be skipped. We note those with an asterisk in the section title, such as *Learning More**.

1.5.2 About the Data

Most of the data sets that we analyze in this book are *simulated* data sets. They are created with R code to have a specific structure. This has several advantages:

- It allows us to illustrate analyses where there is no publicly available marketing data. This is valuable because few firms share their proprietary data for analyses such as segmentation.
- It allows the book to be more self-contained and less dependent on data downloads.
- It makes it possible to alter the data and rerun analyses to see how the results change.
- It lets us teach important R skills for handling data, generating random numbers, and looping in code.

- It demonstrates how one can write analysis code while waiting for real data. When the final data arrives, you can run your code on the new data.

An exception to this is the transactional data in Chap. 12; such data is complex to create and appropriate data has been published [20].

We recommend to work through data simulation sections where they appear; they are designed to teach R and to illustrate points that are typical of marketing data. However, when you need data quickly to continue with a chapter, it is available for download as noted in the next section and again in each chapter.

Whenever possible you should also try to perform the analyses here with your own data sets. We work with data in every chapter, but the best way to learn is to adapt the analyses to other data and work through the issues that arise. Because this is an educational text, not a cookbook, and because R can be slow going at first, we recommend to conduct such parallel analyses on tasks where you are not facing urgent deadlines.

At the beginning, it may seem overly simple to repeat analyses with your own data, but when you try to apply an advanced model to another data set, you'll be much better prepared if you've practiced with multiple data sets all along. The sooner you apply R to your own data, the sooner you will be productive in R.

1.5.3 Online Material

This book has a companion website: <http://r-marketing.r-forge.r-project.org>. The website exists primarily to host the R code and data sets for download, although we encourage you to use those sparingly; you'll learn more if you type the code and create the data sets by simulation as we describe.

On the website, you'll find:

- A welcome page for news and updates: <http://r-marketing.r-forge.r-project.org>
- Code files in `.R` (text) format: <http://r-marketing.r-forge.r-project.org/code>
- Copies of data sets that are simulated in the book: <http://r-marketing.r-forge.r-project.org/data>. These can be downloaded directly into R using the `read.csv()` command (you'll see that command in Sect. 2.6.2, and will find code for an example download in Sect. 3.1)
- A ZIP file containing all of the data and code files: <http://r-marketing.r-forge.r-project.org/data/chapman-feit-rintro.zip>

Links to online data are provided in the form of shortened `goo.gl` links to save typing. More detail on the online materials and ways to access the data are described in Appendix D.

1.5.4 When Things Go Wrong

When you learn something as complex as R or new statistical models, you will encounter many large and small warnings and errors. Also, the R ecosystem is dynamic and things will change after this book is published. We don't wish to scare you with a list of concerns, but we do want you to feel reassured about small discrepancies and to know what to do when larger bugs arise. Here are a few things to know and to try if one of your results doesn't match this book:

- **With R.** The basic error correction process when working with R is to check everything very carefully, especially parentheses, brackets, and upper- or lowercase letters. If a command is lengthy, deconstruct it into pieces and build it up again (we show examples of this along the way).
- **With packages** (add-on libraries). Packages are regularly updated. Sometimes they change how they work, or may not work at all for a while. Some are very stable while others change often. If you have trouble installing one, do a web search for the error message. If output or details are slightly different than we show, don't worry about it. The error "There is no package called . . ." indicates that you need to install the package (Sect. 2.2). For other problems, see the remaining items here or check the package's help file (Sect. 2.4.2).
- **With R warnings and errors.** An R "warning" is often informational and does not necessarily require correction. We call these out as they occur with our code, although sometimes they come and go as packages are updated. If R gives you an "error," that means something went wrong and needs to be corrected. In that case, try the code again, or search online for the error message.
- **With data.** Our data sets are simulated and are affected by random number sequences. If you generate data and it is slightly different, try it again from the beginning; or load the data from the book's website (Sect. 1.5.3).
- **With models.** There are three things that might cause statistical estimates to vary: slight differences in the data (see the preceding item), changes in a package that lead to slightly different estimates, and statistical models that employ random sampling. If you run a model and the results are very similar but slightly different, you can assume that one of these situations occurred. Just proceed.
- **With output.** Packages sometimes change the information they report. The output in this book was current at the time of writing, but you can expect some packages will report things slightly differently over time.
- **With names that can't be located.** Sometimes packages change the function names they use or the structure of results. If you get a code error when trying to extract something from a statistical model, check the model's help file (Sect. 2.4.2); it may be that something has changed names.

Our overall recommendation is this. If the difference is small—such as the difference between a mean of 2.08 and 2.076, or a p -value of 0.726 vs. 0.758—don't

worry too much about it; you can usually safely ignore these. If you find a large difference—such as a statistical estimate of 0.56 instead of 31.92—try the code block again in the book’s code file (Sect. [1.5.3](#)).

1.6 Key Points

At the end of each chapter we summarize crucial lessons. For this chapter, there is only one key point: if you’re ready to learn R, let’s get started with [Chap. 2](#)!

An Overview of the R Language

2.1 Getting Started

In this chapter, we cover just enough of the R language to get you going. If you're new to programming, this chapter will get you started well enough to be productive and we'll call out ways to learn more at the end. R is a great place to learn to program because its environment is clean and much simpler than traditional programming languages such as Java or C++. If you're an experienced programmer in another language, you should skim this chapter to learn the essentials.

We recommend you work through this chapter *hands-on* and be patient; it will prepare you for marketing analytics applications in later chapters.

2.1.1 Initial Steps

If you haven't already installed R, please do so. We'll skip the installation details except to say that you'll want at least the basic version of R (known as "R base") from the comprehensive R archive network (CRAN): <http://cran.r-project.org>. If you are using:

- **Windows or Mac OS X:** Get the *compiled binary* version from CRAN.
- **Linux:** Use your package installer to add R. This might be a GUI installer as in Ubuntu's Software Center or a terminal command such as `sudo apt-get install R`. (See CRAN for more options.)

In either case, you don't need the *source code* version for purposes of this book.

After installing R, we recommend you also install RStudio [140], an integrated environment for writing R code, viewing plots, and reading documentation. RStudio is available for Windows, Mac OS X, and Linux at <http://www.rstudio.com>.

Most users will want the *desktop* version. RStudio is optional and this book does not assume that you're using it, although many R users find it to be convenient. Some companies may have questions about RStudio's Affero general public license (AGPL) terms; if relevant, ask your technology support group if they allow AGPL open source software.

There are other variants of R available, including options that will appeal to experienced programmers who use Emacs, Eclipse, or other development environments. For more information on various R environments, see [Appendix A](#).

2.1.2 Starting R

Once R is installed, run it; or if you installed RStudio, launch that. The R command line starts by default and is known as the *R console*. When this book was written, the R console looked like [Fig. 2.1](#) (where some details depend on the version and operating system).

```
R version 3.1.1 (2014-07-10) -- "Sock it to Me"
Copyright (C) 2014 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin13.1.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[R.app GUI 1.65 (6784) x86_64-apple-darwin13.1.0]
>
```

Fig. 2.1. The R console.

The “>” symbol at the bottom of the R console shows that R is ready for input from you. For example, you could type:

```
> x <- c(2, 4, 6, 8)
```