

## Praise for the First Edition

*Lucid and engaging—this is without doubt the fun way to learn R!*

—Amos A. Folarin, University College London

*Be prepared to quickly raise the bar with the sheer quality that R can produce.*

—Patrick Breen, Rogers Communications Inc.

*An excellent introduction and reference on R from the author of the best R website.*

—Christopher Williams, University of Idaho

*Thorough and readable. A great R companion for the student or researcher.*

—Samuel McQuillin, University of South Carolina

*Finally, a comprehensive introduction to R for programmers.*

—Philipp K. Janert, Author of *Gnuplot in Action*

*Essential reading for anybody moving to R for the first time.*

—Charles Malpas, University of Melbourne

*One of the quickest routes to R proficiency. You can buy the book on Friday and have a working program by Monday.*

—Elizabeth Ostrowski, Baylor College of Medicine

*One usually buys a book to solve the problems they know they have. This book solves problems you didn't know you had.*

—Carles Fenollosa, Barcelona Supercomputing Center

*Clear, precise, and comes with a lot of explanations and examples...the book can be used by beginners and professionals alike, and even for teaching R!*

—Atef Ouni, Tunisian National Institute of Statistics

*A great balance of targeted tutorials and in-depth examples.*

—Landon Cox, 360VL Inc.



# *R in Action*

SECOND EDITION

*Data analysis and graphics with R*

ROBERT I. KABACOFF



MANNING  
SHELTER ISLAND

For online information and ordering of this and other Manning books, please visit [www.manning.com](http://www.manning.com). The publisher offers discounts on this book when ordered in quantity. For more information, please contact


Special Sales Department  
Manning Publications Co.  
20 Baldwin Road  
PO Box 761  
Shelter Island, NY 11964  
Email: [orders@manning.com](mailto:orders@manning.com)

©2015 by Manning Publications Co. All rights reserved.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by means electronic, mechanical, photocopying, or otherwise, without prior written permission of the publisher.

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

- © Recognizing the importance of preserving what has been written, it is Manning's policy to have the books we publish printed on acid-free paper, and we exert our best efforts to that end. Recognizing also our responsibility to conserve the resources of our planet, Manning books are printed on paper that is at least 15 percent recycled and processed without elemental chlorine.

 Manning Publications Co.  
20 Baldwin Road  
PO Box 761  
Shelter Island, NY 11964

Development editor: Jennifer Stout  
Copyeditor: Tiffany Taylor  
Proofreader: Toma Mulligan  
Typesetter: Marija Tudor  
Cover designer: Marija Tudor

ISBN: 9781617291388

Printed in the United States of America

1 2 3 4 5 6 7 8 9 10 – EBM – 20 19 18 17 16 15

# *brief contents*

---

<b>PART 1</b>	<b>GETTING STARTED .....</b>	<b>1</b>
	1 ▪ Introduction to R	3
	2 ▪ Creating a dataset	20
	3 ▪ Getting started with graphs	46
	4 ▪ Basic data management	71
	5 ▪ Advanced data management	89
<b>PART 2</b>	<b>BASIC METHODS .....</b>	<b>115</b>
	6 ▪ Basic graphs	117
	7 ▪ Basic statistics	137
<b>PART 3</b>	<b>INTERMEDIATE METHODS .....</b>	<b>165</b>
	8 ▪ Regression	167
	9 ▪ Analysis of variance	212
	10 ▪ Power analysis	239
	11 ▪ Intermediate graphs	255
	12 ▪ Resampling statistics and bootstrapping	279

<b>PART 4</b>	<b>ADVANCED METHODS .....</b>	<b>299</b>
13	▪ Generalized linear models	301
14	▪ Principal components and factor analysis	319
15	▪ Time series	340
16	▪ Cluster analysis	369
17	▪ Classification	389
18	▪ Advanced methods for missing data	414
<b>PART 5</b>	<b>EXPANDING YOUR SKILLS .....</b>	<b>435</b>
19	▪ Advanced graphics with ggplot2	437
20	▪ Advanced programming	463
21	▪ Creating a package	491
22	▪ Creating dynamic reports	513
23	▪ Advanced graphics with the lattice package	<i>online only</i>

# contents

---

*preface* xvii  
*acknowledgments* xix  
*about this book* xxi  
*about the cover illustration* xxvii

## PART 1 GETTING STARTED ..... 1

### **1** *Introduction to R* 3

- 1.1 Why use R? 5
- 1.2 Obtaining and installing R 7
- 1.3 Working with R 7
  - Getting started* 8 ▪ *Getting help* 10 ▪ *The workspace* 11
  - Input and output* 13
- 1.4 Packages 15
  - What are packages?* 15 ▪ *Installing a package* 15
  - Loading a package* 15 ▪ *Learning about a package* 16
- 1.5 Batch processing 16
- 1.6 Using output as input: reusing results 17
- 1.7 Working with large datasets 17

- 1.8 Working through an example 18
- 1.9 Summary 19

## 2 *Creating a dataset* 20

- 2.1 Understanding datasets 21
- 2.2 Data structures 22
  - Vectors* 22 ▪ *Matrices* 23 ▪ *Arrays* 24 ▪ *Data frames* 25
  - Factors* 28 ▪ *Lists* 30
- 2.3 Data input 32
  - Entering data from the keyboard* 33 ▪ *Importing data from a delimited text file* 34 ▪ *Importing data from Excel* 37
  - Importing data from XML* 38 ▪ *Importing data from the web* 38 ▪ *Importing data from SPSS* 38 ▪ *Importing data from SAS* 39 ▪ *Importing data from Stata* 40 ▪ *Importing data from NetCDF* 40 ▪ *Importing data from HDF5* 40
  - Accessing database management systems (DBMSs)* 40
  - Importing data via Stat/Transfer* 42
- 2.4 Annotating datasets 43
  - Variable labels* 43 ▪ *Value labels* 43
- 2.5 Useful functions for working with data objects 43
- 2.6 Summary 44

## 3 *Getting started with graphs* 46

- 3.1 Working with graphs 47
- 3.2 A simple example 49
- 3.3 Graphical parameters 50
  - Symbols and lines* 51 ▪ *Colors* 52 ▪ *Text characteristics* 53
  - Graph and margin dimensions* 54
- 3.4 Adding text, customized axes, and legends 56
  - Titles* 56 ▪ *Axes* 57 ▪ *Reference lines* 60 ▪ *Legend* 60
  - Text annotations* 61 ▪ *Math annotations* 63
- 3.5 Combining graphs 64
  - Creating a figure arrangement with fine control* 68
- 3.6 Summary 70

## 4 *Basic data management* 71

- 4.1 A working example 71
- 4.2 Creating new variables 73



- 4.3 Recoding variables 75
  - 4.4 Renaming variables 76
  - 4.5 Missing values 77
    - Recoding values to missing* 78 ▪ *Excluding missing values from analyses* 78
  - 4.6 Date values 79
    - Converting dates to character variables* 81 ▪ *Going further* 81
  - 4.7 Type conversions 81
  - 4.8 Sorting data 82
  - 4.9 Merging datasets 83
    - Adding columns to a data frame* 83 ▪ *Adding rows to a data frame* 84
  - 4.10 Subsetting datasets 84
    - Selecting (keeping) variables* 84 ▪ *Excluding (dropping) variables* 84 ▪ *Selecting observations* 85 ▪ *The subset() function* 86 ▪ *Random samples* 87
  - 4.11 Using SQL statements to manipulate data frames 87
  - 4.12 Summary 88
- 5 *Advanced data management* 89**
- 5.1 A data-management challenge 90
  - 5.2 Numerical and character functions 91
    - Mathematical functions* 91 ▪ *Statistical functions* 92
    - Probability functions* 94 ▪ *Character functions* 97
    - Other useful functions* 98 ▪ *Applying functions to matrices and data frames* 99
  - 5.3 A solution for the data-management challenge 101
  - 5.4 Control flow 105
    - Repetition and looping* 105 ▪ *Conditional execution* 106
  - 5.5 User-written functions 107
  - 5.6 Aggregation and reshaping 109
    - Transpose* 110 ▪ *Aggregating data* 110 ▪ *The reshape2 package* 111
  - 5.7 Summary 113

## PART 2 BASIC METHODS ..... 115

### 6 **Basic graphs 117**

- 6.1 Bar plots 118
  - Simple bar plots 118* ▪ *Stacked and grouped bar plots 119*
  - Mean bar plots 120* ▪ *Tweaking bar plots 121*
  - Spinograms 122*
- 6.2 Pie charts 123
- 6.3 Histograms 125
- 6.4 Kernel density plots 127
- 6.5 Box plots 129
  - Using parallel box plots to compare groups 129* ▪ *Violin plots 132*
- 6.6 Dot plots 133
- 6.7 Summary 136

### 7 **Basic statistics 137**

- 7.1 Descriptive statistics 138
  - A menagerie of methods 138* ▪ *Even more methods 140*
  - Descriptive statistics by group 142* ▪ *Additional methods by group 143* ▪ *Visualizing results 144*
- 7.2 Frequency and contingency tables 144
  - Generating frequency tables 145* ▪ *Tests of independence 151* ▪ *Measures of association 152*
  - Visualizing results 153*
- 7.3 Correlations 153
  - Types of correlations 153* ▪ *Testing correlations for significance 156* ▪ *Visualizing correlations 158*
- 7.4 T-tests 158
  - Independent t-test 158* ▪ *Dependent t-test 159*
  - When there are more than two groups 160*
- 7.5 Nonparametric tests of group differences 160
  - Comparing two groups 160* ▪ *Comparing more than two groups 161*
- 7.6 Visualizing group differences 163
- 7.7 Summary 164

## PART 3 INTERMEDIATE METHODS ..... 165

8 **Regression 167**

- 8.1 The many faces of regression 168
  - Scenarios for using OLS regression 169* ▪ *What you need to know 170*
- 8.2 OLS regression 171
  - Fitting regression models with lm() 172* ▪ *Simple linear regression 173* ▪ *Polynomial regression 175*
  - Multiple linear regression 178* ▪ *Multiple linear regression with interactions 180*
- 8.3 Regression diagnostics 182
  - A typical approach 183* ▪ *An enhanced approach 187*
  - Global validation of linear model assumption 193*
  - Multicollinearity 193*
- 8.4 Unusual observations 194
  - Outliers 194* ▪ *High-leverage points 195* ▪ *Influential observations 196*
- 8.5 Corrective measures 198
  - Deleting observations 199* ▪ *Transforming variables 199*
  - Adding or deleting variables 201* ▪ *Trying a different approach 201*
- 8.6 Selecting the “best” regression model 201
  - Comparing models 202* ▪ *Variable selection 203*
- 8.7 Taking the analysis further 206
  - Cross-validation 206* ▪ *Relative importance 208*
- 8.8 Summary 211

9 **Analysis of variance 212**

- 9.1 A crash course on terminology 213
- 9.2 Fitting ANOVA models 215
  - The aov() function 215* ▪ *The order of formula terms 216*
- 9.3 One-way ANOVA 218
  - Multiple comparisons 219* ▪ *Assessing test assumptions 222*
- 9.4 One-way ANCOVA 223
  - Assessing test assumptions 225* ▪ *Visualizing the results 225*
- 9.5 Two-way factorial ANOVA 226

- 9.6 Repeated measures ANOVA 229
- 9.7 Multivariate analysis of variance (MANOVA) 232
  - Assessing test assumptions* 234 ▪ *Robust MANOVA* 235
- 9.8 ANOVA as regression 236
- 9.9 Summary 238

## 10 *Power analysis* 239

- 10.1 A quick review of hypothesis testing 240
- 10.2 Implementing power analysis with the pwr package 242
  - t-tests* 243 ▪ *ANOVA* 245 ▪ *Correlations* 245
  - Linear models* 246 ▪ *Tests of proportions* 247
  - Chi-square tests* 248 ▪ *Choosing an appropriate effect size in novel situations* 249
- 10.3 Creating power analysis plots 251
- 10.4 Other packages 252
- 10.5 Summary 253

## 11 *Intermediate graphs* 255

- 11.1 Scatter plots 256
  - Scatter-plot matrices* 259 ▪ *High-density scatter plots* 261
  - 3D scatter plots* 263 ▪ *Spinning 3D scatter plots* 265
  - Bubble plots* 266
- 11.2 Line charts 268
- 11.3 Corrgrams 271
- 11.4 Mosaic plots 276
- 11.5 Summary 278

## 12 *Resampling statistics and bootstrapping* 279

- 12.1 Permutation tests 280
- 12.2 Permutation tests with the coin package 282
  - Independent two-sample and k-sample tests* 283
  - Independence in contingency tables* 285 ▪ *Independence between numeric variables* 285 ▪ *Dependent two-sample and k-sample tests* 286 ▪ *Going further* 286
- 12.3 Permutation tests with the lmPerm package 287
  - Simple and polynomial regression* 287 ▪ *Multiple regression* 288 ▪ *One-way ANOVA and ANCOVA* 289
  - Two-way ANOVA* 290

- 12.4 Additional comments on permutation tests 291
- 12.5 Bootstrapping 291
- 12.6 Bootstrapping with the boot package 292
  - Bootstrapping a single statistic* 294 ▪ *Bootstrapping several statistics* 296
- 12.7 Summary 298

## PART 4 ADVANCED METHODS ..... 299

### 13 *Generalized linear models* 301

- 13.1 Generalized linear models and the glm() function 302
  - The glm() function* 303 ▪ *Supporting functions* 304
  - Model fit and regression diagnostics* 305
- 13.2 Logistic regression 306
  - Interpreting the model parameters* 308 ▪ *Assessing the impact of predictors on the probability of an outcome* 309
  - Overdispersion* 310 ▪ *Extensions* 311
- 13.3 Poisson regression 312
  - Interpreting the model parameters* 314 ▪ *Overdispersion* 315
  - Extensions* 317
- 13.4 Summary 318

### 14 *Principal components and factor analysis* 319

- 14.1 Principal components and factor analysis in R 321
- 14.2 Principal components 322
  - Selecting the number of components to extract* 323
  - Extracting principal components* 324 ▪ *Rotating principal components* 327 ▪ *Obtaining principal components scores* 328
- 14.3 Exploratory factor analysis 330
  - Deciding how many common factors to extract* 331
  - Extracting common factors* 332 ▪ *Rotating factors* 333
  - Factor scores* 336 ▪ *Other EFA-related packages* 337
- 14.4 Other latent variable models 337
- 14.5 Summary 338

### 15 *Time series* 340

- 15.1 Creating a time-series object in R 343

- 15.2 Smoothing and seasonal decomposition 345
  - Smoothing with simple moving averages* 345
  - *Seasonal decomposition* 347
- 15.3 Exponential forecasting models 352
  - Simple exponential smoothing* 353
  - *Holt and Holt-Winters exponential smoothing* 355
  - *The ets() function and automated forecasting* 358
- 15.4 ARIMA forecasting models 359
  - Prerequisite concepts* 359
  - *ARMA and ARIMA models* 361
  - Automated ARIMA forecasting* 366
- 15.5 Going further 367
- 15.6 Summary 367

## 16 *Cluster analysis* 369

- 16.1 Common steps in cluster analysis 370
- 16.2 Calculating distances 372
- 16.3 Hierarchical cluster analysis 374
- 16.4 Partitioning cluster analysis 378
  - K-means clustering* 378
  - *Partitioning around medoids* 382
- 16.5 Avoiding nonexistent clusters 384
- 16.6 Summary 387

## 17 *Classification* 389

- 17.1 Preparing the data 390
- 17.2 Logistic regression 392
- 17.3 Decision trees 393
  - Classical decision trees* 393
  - *Conditional inference trees* 397
- 17.4 Random forests 399
- 17.5 Support vector machines 401
  - Tuning an SVM* 403
- 17.6 Choosing a best predictive solution 405
- 17.7 Using the rattle package for data mining 408
- 17.8 Summary 413

## 18 *Advanced methods for missing data* 414

- 18.1 Steps in dealing with missing data 415
- 18.2 Identifying missing values 417

- 18.3 Exploring missing-values patterns 418
  - Tabulating missing values* 419 ▪ *Exploring missing data visually* 419 ▪ *Using correlations to explore missing values* 422
- 18.4 Understanding the sources and impact of missing data 424
- 18.5 Rational approaches for dealing with incomplete data 425
- 18.6 Complete-case analysis (listwise deletion) 426
- 18.7 Multiple imputation 428
- 18.8 Other approaches to missing data 432
  - Pairwise deletion* 432 ▪ *Simple (nonstochastic) imputation* 433
- 18.9 Summary 433

## PART 5 EXPANDING YOUR SKILLS ..... 435

### 19 *Advanced graphics with ggplot2* 437

- 19.1 The four graphics systems in R 438
- 19.2 An introduction to the ggplot2 package 439
- 19.3 Specifying the plot type with geoms 443
- 19.4 Grouping 447
- 19.5 Faceting 450
- 19.6 Adding smoothed lines 453
- 19.7 Modifying the appearance of ggplot2 graphs 455
  - Axes* 455 ▪ *Legends* 457 ▪ *Scales* 458 ▪ *Themes* 460
  - Multiple graphs per page* 461
- 19.8 Saving graphs 462
- 19.9 Summary 462

### 20 *Advanced programming* 463

- 20.1 A review of the language 464
  - Data types* 464 ▪ *Control structures* 470 ▪ *Creating functions* 473
- 20.2 Working with environments 475
- 20.3 Object-oriented programming 477
  - Generic functions* 477 ▪ *Limitations of the S3 model* 479
- 20.4 Writing efficient code 479

- 20.5 Debugging 483
  - Common sources of errors* 483 ▪ *Debugging tools* 484
  - Session options that support debugging* 486
- 20.6 Going further 489
- 20.7 Summary 490

## 21 *Creating a package* 491

- 21.1 Nonparametric analysis and the *npar* package 492
  - Comparing groups with the npar package* 494
- 21.2 Developing the package 496
  - Computing the statistics* 497 ▪ *Printing the results* 500
  - Summarizing the results* 501 ▪ *Plotting the results* 504
  - Adding sample data to the package* 505
- 21.3 Creating the package documentation 506
- 21.4 Building the package 508
- 21.5 Going further 512
- 21.6 Summary 512

## 22 *Creating dynamic reports* 513

- 22.1 A template approach to reports 515
- 22.2 Creating dynamic reports with R and Markdown 517
- 22.3 Creating dynamic reports with R and LaTeX 522
- 22.4 Creating dynamic reports with R and Open Document 525
- 22.5 Creating dynamic reports with R and Microsoft Word 527
- 22.6 Summary 531

- afterword* *Into the rabbit hole* 532
- appendix A* *Graphical user interfaces* 535
- appendix B* *Customizing the startup environment* 538
- appendix C* *Exporting data from R* 540
- appendix D* *Matrix algebra in R* 542
- appendix E* *Packages used in this book* 544
- appendix F* *Working with large datasets* 551
- appendix G* *Updating an R installation* 555
- references* 558
- index* 563

- bonus chapter 23* *Advanced graphics with the lattice package*  
available online at [manning.com/RinActionSecondEdition](http://manning.com/RinActionSecondEdition)



# *preface*

---

*What is the use of a book, without pictures or conversations?*

—Alice, *Alice's Adventures in Wonderland*

*It's wondrous, with treasures to satiate desires both subtle and gross; but it's not for the timid.*

—Q, "Q Who?" *Stark Trek: The Next Generation*

When I began writing this book, I spent quite a bit of time searching for a good quote to start things off. I ended up with two. R is a wonderfully flexible platform and language for exploring, visualizing, and understanding data. I chose the quote from *Alice's Adventures in Wonderland* to capture the flavor of statistical analysis today—an interactive process of exploration, visualization, and interpretation.

The second quote reflects the generally held notion that R is difficult to learn. What I hope to show you is that it doesn't have to be. R is broad and powerful, with so many analytic and graphic functions available (more than 50,000 at last count) that it easily intimidates both novice and experienced users alike. But there is rhyme and reason to the apparent madness. With guidelines and instructions, you can navigate the tremendous resources available, selecting the tools you need to accomplish your work with style, elegance, efficiency—and more than a little coolness.

I first encountered R several years ago, when applying for a new statistical consulting position. The prospective employer asked in the pre-interview material if I was conversant in R. Following the standard advice of recruiters, I immediately said yes,

and set off to learn it. I was an experienced statistician and researcher, had 25 years experience as an SAS and SPSS programmer, and was fluent in a half dozen programming languages. How hard could it be? Famous last words.

As I tried to learn the language (as fast as possible, with an interview looming), I found either tomes on the underlying structure of the language or dense treatises on specific advanced statistical methods, written by and for subject-matter experts. The online help was written in a spartan style that was more reference than tutorial. Every time I thought I had a handle on the overall organization and capabilities of R, I found something new that made me feel ignorant and small.

To make sense of it all, I approached R as a data scientist. I thought about what it takes to successfully process, analyze, and understand data, including

- Accessing the data (getting the data into the application from multiple sources)
- Cleaning the data (coding missing data, fixing or deleting miscoded data, transforming variables into more useful formats)
- Annotating the data (in order to remember what each piece represents)
- Summarizing the data (getting descriptive statistics to help characterize the data)
- Visualizing the data (because a picture really is worth a thousand words)
- Modeling the data (uncovering relationships and testing hypotheses)
- Preparing the results (creating publication-quality tables and graphs)

Then I tried to understand how I could use R to accomplish each of these tasks. Because I learn best by teaching, I eventually created a website ([www.statmethods.net](http://www.statmethods.net)) to document what I had learned.

Then, about a year later, Marjan Bace, Manning's publisher, called and asked if I would like to write a book on R. I had already written 50 journal articles, 4 technical manuals, numerous book chapters, and a book on research methodology, so how hard could it be? At the risk of sounding repetitive—famous last words.

A year after the first edition came out in 2011, I started working on the second edition. The R platform is evolving, and I wanted to describe these new developments. I also wanted to expand the coverage of predictive analytics and data mining—important topics in the world of big data. Finally, I wanted to add chapters on advanced data visualization, software development, and dynamic report writing.

The book you're holding is the one that I wished I had so many years ago. I have tried to provide you with a guide to R that will allow you to quickly access the power of this great open source endeavor, without all the frustration and angst. I hope you enjoy it.

P.S. I was offered the job but didn't take it. But learning R has taken my career in directions that I could never have anticipated. Life can be funny.

# *acknowledgments*

---

A number of people worked hard to make this a better book. They include

- Marjan Bace, Manning's publisher, who asked me to write this book in the first place.
- Sebastian Stirling and Jennifer Stout, development editors on the first and second editions, respectively. Each spent many hours helping me organize the material, clarify concepts, and generally make the text more interesting.
- Pablo Domínguez Vaselli, technical proofreader, who helped uncover areas of confusion and provided an independent and expert eye for testing code. I came to rely on his vast knowledge, careful reviews, and considered judgment.
- Olivia Booth, the review editor, who helped obtain reviewers and coordinate the review process.
- Mary Piergies, who helped shepherd this book through the production process, and her team of Tiffany Taylor, Toma Mulligan, Janet Vail, David Novak, and Marija Tudor.
- The peer reviewers who spent hours of their own time carefully reading through the material, finding typos, and making valuable substantive suggestions: Bryce Darling, Christian Theil Have, Cris Weber, Deepak Vohra, Dwight Barry, George Gaines, Indrajit Sen Gupta, Dr. L. Duleep Kumar Samuel, Mahesh Srinivason, Marc Paradis, Peter Rabinovitch, Ravishankar Rajagopalan, Samuel Dale McQuillin, and Zekai Otles.
- The many Manning Early Access Program (MEAP) participants who bought the book before it was finished, asked great questions, pointed out errors, and made helpful suggestions.

Each contributor has made this a better and more comprehensive book.

I would also like to acknowledge the many software authors who have contributed to making R such a powerful data-analytic platform. They include not only the core developers, but also the selfless individuals who have created and maintain contributed packages, extending R's capabilities greatly. Appendix E provides a list of the authors of contributed packages described in this book. In particular, I would like to mention John Fox, Hadley Wickham, Frank E. Harrell, Jr., Deepayan Sarkar, and William Revelle, whose works I greatly admire. I have tried to represent their contributions accurately, and I remain solely responsible for any errors or distortions inadvertently included in this book.

I really should have started this book by thanking my wife and partner, Carol Lynn. Although she has no intrinsic interest in statistics or programming, she read each chapter multiple times and made countless corrections and suggestions. No greater love has any person than to read multivariate statistics for another. Just as important, she suffered the long nights and weekends that I spent writing this book, with grace, support, and affection. There is no logical explanation why I should be this lucky.

There are two other people I would like to thank. One is my father, whose love of science was inspiring and who gave me an appreciation of the value of data. I miss him dearly. The other is Gary K. Burger, my mentor in graduate school. Gary got me interested in a career in statistics and teaching when I thought I wanted to be a clinician. This is all his fault.

## *about this book*

---

If you picked up this book, you probably have some data that you need to collect, summarize, transform, explore, model, visualize, or present. If so, then R is for you! R has become the worldwide language for statistics, predictive analytics, and data visualization. It offers the widest range of methodologies for understanding data currently available, from the most basic to the most complex and bleeding edge.

As an open source project it's freely available for a range of platforms, including Windows, Mac OS X, and Linux. It's under constant development, with new procedures added daily. Additionally, R is supported by a large and diverse community of data scientists and programmers who gladly offer their help and advice to users.

Although R is probably best known for its ability to create beautiful and sophisticated graphs, it can handle just about any statistical problem. The base installation provides hundreds of data-management, statistical, and graphical functions out of the box. But some of its most powerful features come from the thousands of extensions (packages) provided by contributing authors.

This breadth comes at a price. It can be hard for new users to get a handle on what R is and what it can do. Even the most experienced R user is surprised to learn about features they were unaware of.

*R in Action, Second Edition* provides you with a guided introduction to R, giving you a 2,000-foot view of the platform and its capabilities. It will introduce you to the most important functions in the base installation and more than 90 of the most useful contributed packages. Throughout the book, the goal is practical application—how you can make sense of your data and communicate that understanding to others. When you finish, you should have a good grasp of how R works and what it can do and where

you can go to learn more. You'll be able to apply a variety of techniques for visualizing data, and you'll have the skills to tackle both basic and advanced data analytic problems.

### ***What's new in the second edition***

If you want to delve into the use of R more deeply, the second edition offers more than 200 pages of new material. Concentrated in the second half of the book are new chapters on data mining, predictive analytics, and advanced programming. In particular, chapters 15 (time series), 16 (cluster analysis), 17 (classification), 19 (ggplot2 graphics), 20 (advanced programming), 21 (creating a package), and 22 (creating dynamic reports) are new. In addition, chapter 2 (creating a dataset) has more detailed information on importing data from text and SAS files, and appendix F (working with large datasets) has been expanded to include new tools for working with big data problems. Finally, numerous updates and corrections have been made throughout the text.

### ***Who should read this book***

*R in Action, Second Edition* should appeal to anyone who deals with data. No background in statistical programming or the R language is assumed. Although the book is accessible to novices, there should be enough new and practical material to satisfy even experienced R mavens.

Users without a statistical background who want to use R to manipulate, summarize, and graph data should find chapters 1–6, 11, and 19 easily accessible. Chapters 7 and 10 assume a one-semester course in statistics; and readers of chapters 8, 9, and 12–18 will benefit from two semesters of statistics. Chapters 20–22 offer a deeper dive into the R language and have no statistical prerequisites. I've tried to write each chapter in such a way that both beginning and expert data analysts will find something interesting and useful.

### ***Roadmap***

This book is designed to give you a guided tour of the R platform, with a focus on those methods most immediately applicable for manipulating, visualizing, and understanding data. The book has 22 chapters and is divided into 5 parts: "Getting Started," "Basic Methods," "Intermediate Methods," "Advanced Methods," and "Expanding Your Skills." Additional topics are covered in seven appendices.

Chapter 1 begins with an introduction to R and the features that make it so useful as a data-analysis platform. The chapter covers how to obtain the program and how to enhance the basic installation with extensions that are available online. The remainder of the chapter is spent exploring the user interface and learning how to run programs interactively and in batch.

Chapter 2 covers the many methods available for getting data into R. The first half of the chapter introduces the data structures R uses to hold data, and how to enter

data from the keyboard. The second half discusses methods for importing data into R from text files, web pages, spreadsheets, statistical packages, and databases.

Many users initially approach R because they want to create graphs, so we jump right into that topic in chapter 3. No waiting required. We review methods of creating graphs, modifying them, and saving them in a variety of formats.

Chapter 4 covers basic data management, including sorting, merging, and subsetting datasets, and transforming, recoding, and deleting variables.

Building on the material in chapter 4, chapter 5 covers the use of functions (mathematical, statistical, character) and control structures (looping, conditional execution) for data management. I then discuss how to write your own R functions and how to aggregate data in various ways.

Chapter 6 demonstrates methods for creating common univariate graphs, such as bar plots, pie charts, histograms, density plots, box plots, and dot plots. Each is useful for understanding the distribution of a single variable.

Chapter 7 starts by showing how to summarize data, including the use of descriptive statistics and cross-tabulations. We then look at basic methods for understanding relationships between two variables, including correlations, t-tests, chi-square tests, and nonparametric methods.

Chapter 8 introduces regression methods for modeling the relationship between a numeric outcome variable and a set of one or more numeric predictor variables. Methods for fitting these models, evaluating their appropriateness, and interpreting their meaning are discussed in detail.

Chapter 9 considers the analysis of basic experimental designs through the analysis of variance and its variants. Here we're usually interested in how treatment combinations or conditions affect a numerical outcome. Methods for assessing the appropriateness of the analyses and visualizing the results are also covered.

Chapter 10 provides a detailed treatment of power analysis. Starting with a discussion of hypothesis testing, the chapter focuses on how to determine the sample size necessary to detect a treatment effect of a given size with a given degree of confidence. This can help you to plan experimental and quasi-experimental studies that are likely to yield useful results.

Chapter 11 expands on the material in chapter 6, covering the creation of graphs that help you to visualize relationships among two or more variables. These include various types of 2D and 3D scatter plots, scatter-plot matrices, line plots, correlograms, and mosaic plots.

Chapter 12 presents analytic methods that work well in cases where data are sampled from unknown or mixed distributions, where sample sizes are small, where outliers are a problem, or where devising an appropriate test based on a theoretical distribution is too complex and mathematically intractable. They include both resampling and bootstrapping approaches—computer-intensive methods that are easily implemented in R.

Chapter 13 expands on the regression methods in chapter 8 to cover data that are not normally distributed. The chapter starts with a discussion of generalized linear

models and then focuses on cases where you're trying to predict an outcome variable that is either categorical (logistic regression) or a count (Poisson regression).

One of the challenges of multivariate data problems is simplification. Chapter 14 describes methods of transforming a large number of correlated variables into a smaller set of uncorrelated variables (principal component analysis), as well as methods for uncovering the latent structure underlying a given set of variables (factor analysis). The many steps involved in an appropriate analysis are covered in detail.

Chapter 15 describes methods for creating, manipulating, and modeling time series data. It covers visualizing and decomposing time series data, as well as exponential and ARIMA approaches to forecasting future values.

Chapter 16 illustrates methods of clustering observations into naturally occurring groups. The chapter begins with a discussion of the common steps in a comprehensive cluster analysis, followed by a presentation of hierarchical clustering and partitioning methods. Several methods for determining the proper number of clusters are presented.

Chapter 17 presents popular supervised machine-learning methods for classifying observations into groups. Decision trees, random forests, and support vector machines are considered in turn. You'll also learn about methods for evaluating the accuracy of each approach.

In keeping with my attempt to present practical methods for analyzing data, chapter 18 considers modern approaches to the ubiquitous problem of missing data values. R supports a number of elegant approaches for analyzing datasets that are incomplete for various reasons. Several of the best are described here, along with guidance for which ones to use when, and which ones to avoid.

Chapter 19 wraps up the discussion of graphics with a presentation of one of R's most useful and advanced approaches to visualizing data: `ggplot2`. The `ggplot2` package implements a grammar of graphics that provides a powerful and consistent set of tools for graphing multivariate data.

Chapter 20 covers advanced programming techniques. You'll learn about object-oriented programming techniques and debugging approaches. The chapter also presents a variety of tips for efficient programming. This chapter will be particularly helpful if you're seeking a greater understanding of how R works, and it's a prerequisite for chapter 21.

Chapter 21 provides a step-by-step guide to creating R packages. This will allow you to create more sophisticated programs, document them efficiently, and share them with others.

Finally, chapter 22 offers several methods for creating attractive reports from within R. You'll learn how to generate web pages, reports, articles, and even books from your R code. The resulting documents can include your code, tables of results, graphs, and commentary.

The afterword points you to many of the best internet sites for learning more about R, joining the R community, getting questions answered, and staying current with this rapidly changing product.



Last, but not least, the seven appendices (A through G) extend the text's coverage to include such useful topics as R graphic user interfaces, customizing and upgrading an R installation, exporting data to other applications, using R for matrix algebra (à la MATLAB), and working with very large datasets.

We also offer a bonus chapter, which is available online only from the publisher's website at [manning.com/RinActionSecondEdition](http://manning.com/RinActionSecondEdition). Online chapter 23 covers the lattice package, which is introduced in chapter 19.

### **Advice for data miners**

Data mining is a field of analytics concerned with discovering patterns in large data sets. Many data-mining specialists are turning to R for its cutting-edge analytical capabilities. If you're a data miner making the transition to R and want to access the language as quickly as possible, I recommend the following reading sequence: chapter 1 (introduction), chapter 2 (data structures and those portions of importing data that are relevant to your setting), chapter 4 (basic data management), chapter 7 (descriptive statistics), chapter 8 (sections 1, 2, and 6; regression), chapter 13 (section 2; logistic regression), chapter 16 (clustering), chapter 17 (classification), and appendix F (working with large datasets). Then review the other chapters as needed.

### **Code examples**

In order to make this book as broadly applicable as possible, I've chosen examples from a range of disciplines, including psychology, sociology, medicine, biology, business, and engineering. None of these examples require a specialized knowledge of that field.

The datasets used in these examples were selected because they pose interesting questions and because they're small. This allows you to focus on the techniques described and quickly understand the processes involved. When you're learning new methods, smaller is better. The datasets are provided with the base installation of R or available through add-on packages that are available online.

The source code for each example is available from [www.manning.com/RinActionSecondEdition](http://www.manning.com/RinActionSecondEdition) and at [www.github.com/kabacoff/RiA2](http://www.github.com/kabacoff/RiA2). To get the most out of this book, I recommend that you try the examples as you read them.

Finally, a common maxim states that if you ask two statisticians how to analyze a dataset, you'll get three answers. The flip side of this assertion is that each answer will move you closer to an understanding of the data. I make no claim that a given analysis is the best or only approach to a given problem. Using the skills taught in this text, I invite you to play with the data and see what you can learn. R is interactive, and the best way to learn is to experiment.

### **Code conventions**

The following typographical conventions are used throughout this book:

- A monospaced font is used for code listings that should be typed as is.

- A monospaced font is also used within the general text to denote code words or previously defined objects.
- *Italics* within code listings indicate placeholders. You should replace them with appropriate text and values for the problem at hand. For example, `path_to_my_file` would be replaced with the actual path to a file on your computer.
- R is an interactive language that indicates readiness for the next line of user input with a prompt (`>` by default). Many of the listings in this book capture interactive sessions. When you see code lines that start with `>`, don't type the prompt.
- Code annotations are used in place of inline comments (a common convention in Manning books). Additionally, some annotations appear with numbered bullets like ❶ that refer to explanations appearing later in the text.
- To save room or make text more legible, the output from interactive sessions may include additional white space or omit text that is extraneous to the point under discussion.

### **Author Online**

Purchase of *R in Action, Second Edition* includes free access to a private web forum run by Manning Publications where you can make comments about the book, ask technical questions, and receive help from the author and from other users. To access the forum and subscribe to it, point your web browser to [www.manning.com/RinActionSecondEdition](http://www.manning.com/RinActionSecondEdition). This page provides information on how to get on the forum once you're registered, what kind of help is available, and the rules of conduct on the forum.

Manning's commitment to our readers is to provide a venue where a meaningful dialog between individual readers and between readers and the author can take place. It isn't a commitment to any specific amount of participation on the part of the author, whose contribution to the AO forum remains voluntary (and unpaid). We suggest you try asking the author some challenging questions, lest his interest stray!

The AO forum and the archives of previous discussions will be accessible from the publisher's website as long as the book is in print.

### **About the author**

Dr. Robert Kabacoff is Vice President of Research for Management Research Group, an international organizational development and consulting firm. He has more than 20 years of experience providing research and statistical consultation to organizations in health care, financial services, manufacturing, behavioral sciences, government, and academia. Prior to joining MRG, Dr. Kabacoff was a professor of psychology at Nova Southeastern University in Florida, where he taught graduate courses in quantitative methods and statistical programming. For the past five years, he has managed Quick-R ([www.statmethods.net](http://www.statmethods.net)), a popular R tutorial website.

## *about the cover illustration*

---

The figure on the cover of *R in Action, Second Edition* is captioned “A man from Zadar.” The illustration is taken from a reproduction of an album of Croatian traditional costumes from the mid-nineteenth century by Nikola Arsenovic, published by the Ethnographic Museum in Split, Croatia, in 2003. The illustrations were obtained from a helpful librarian at the Ethnographic Museum in Split, itself situated in the Roman core of the medieval center of the town: the ruins of Emperor Diocletian’s retirement palace from around AD 304. The book includes finely colored illustrations of figures from different regions of Croatia, accompanied by descriptions of the costumes and of everyday life.

Zadar is an old Roman-era town on the northern Dalmatian coast of Croatia. It’s over 2,000 years old and served for hundreds of years as an important port on the trading route from Constantinople to the West. Situated on a peninsula framed by small Adriatic islands, the city is picturesque and has become a popular tourist destination with its architectural treasures of Roman ruins, moats, and old stone walls. The figure on the cover wears blue woolen trousers and a white linen shirt, over which he dons a blue vest and jacket trimmed with the colorful embroidery typical for this region. A red woolen belt and cap complete the costume.

Dress codes and lifestyles have changed over the last 200 years, and the diversity by region, so rich at the time, has faded away. It’s now hard to tell apart the inhabitants of different continents, let alone of different hamlets or towns separated by only a few miles. Perhaps we have traded this cultural diversity for a more varied personal life—certainly for a more varied and fast-paced technological life.

Manning celebrates the inventiveness and initiative of the computer business with book covers based on the rich diversity of regional life of two centuries ago, brought back to life by illustrations from old books and collections like this one.