

```
> mydf[[2]]  
[1] 3 2 11 4 0  
> mydf[["Cost"]]  
[1] 3 2 11 4 0  
> mydf[["C"]]  
NULL  
> mydf[["C", exact = FALSE]]  
[1] 3 2 11 4 0
```

# A DATA SCIENTIST'S GUIDE TO ACQUIRING, CLEANING, and MANAGING DATA in R

Samuel E. Buttrey and Lyn R. Whitaker



WILEY



**A Data Scientist's Guide to Acquiring,  
Cleaning, and Managing Data in R**



# **A Data Scientist's Guide to Acquiring, Cleaning, and Managing Data in R**

*Samuel E. Buttrey and Lyn R. Whitaker*  
*Naval Postgraduate School, California, United States*

**WILEY**

This edition first published 2018

© 2018 John Wiley & Sons Ltd

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by law. Advice on how to obtain permission to reuse material from this title is available at <http://www.wiley.com/go/permissions>.

The right of Samuel E. Buttrey and Lyn R. Whitaker to be identified as the authors of this work has been asserted in accordance with law.

*Registered Offices*

John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, USA

John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex, PO19 8SQ, UK

*Editorial Office*

9600 Garsington Road, Oxford, OX4 2DQ, UK

For details of our global editorial offices, customer services, and more information about Wiley products visit us at [www.wiley.com](http://www.wiley.com).

Wiley also publishes its books in a variety of electronic formats and by print-on-demand. Some content that appears in standard print versions of this book may not be available in other formats.

*Limit of Liability/Disclaimer of Warranty*

While the publisher and authors have used their best efforts in preparing this work, they make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives, written sales materials or promotional statements for this work. The fact that an organization, website, or product is referred to in this work as a citation and/or potential source of further information does not mean that the publisher and authors endorse the information or services the organization, website, or product may provide or recommendations it may make. This work is sold with the understanding that the publisher is not engaged in rendering professional services. The advice and strategies contained herein may not be suitable for your situation. You should consult with a specialist where appropriate. Further, readers should be aware that websites listed in this work may have changed or disappeared between when this work was written and when it is read. Neither the publisher nor authors shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

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

Hardback ISBN: 9781119080022

Cover Design: Wiley

Cover Image: © Nongkran\_ch/Gettyimages

Set in 10/12pt WarnockPro by SPi Global, Chennai, India

10 9 8 7 6 5 4 3 2 1

*To Linda and Mike*



## Contents

<b>About the Authors</b>	<i>xv</i>
<b>Preface</b>	<i>xvii</i>
<b>Acknowledgments</b>	<i>xix</i>
<b>About the Companion Website</b>	<i>xxi</i>

<b>1</b>	<b>R</b>	<i>1</i>
1.1	Introduction	<i>1</i>
1.1.1	What Is R?	<i>1</i>
1.1.2	Who Uses R and Why?	<i>2</i>
1.1.3	Acquiring and Installing R	<i>2</i>
1.1.4	Starting and Quitting R	<i>3</i>
1.2	Data	<i>3</i>
1.2.1	Acquiring Data	<i>3</i>
1.2.2	Cleaning Data	<i>4</i>
1.2.3	The Goal of Data Cleaning	<i>4</i>
1.2.4	Making Your Work Reproducible	<i>5</i>
1.3	The Very Basics of R	<i>5</i>
1.3.1	Top Ten Quick Facts You Need to Know about R	<i>5</i>
1.3.2	Vocabulary	<i>8</i>
1.3.3	Calculating and Printing in R	<i>11</i>
1.4	Running an R Session	<i>12</i>
1.4.1	Where Your Data Is Stored	<i>13</i>
1.4.2	Options	<i>13</i>
1.4.3	Scripts	<i>14</i>
1.4.4	R Packages	<i>14</i>
1.4.5	RStudio and Other GUIs	<i>15</i>
1.4.6	Locales and Character Sets	<i>15</i>
1.5	Getting Help	<i>16</i>
1.5.1	At the Command Line	<i>16</i>
1.5.2	The Online Manuals	<i>16</i>
1.5.3	On the Internet	<i>17</i>

1.5.4	Further Reading	17
1.6	How to Use This Book	17
1.6.1	Syntax and Conventions in This Book	17
1.6.2	The Chapters	18
<b>2</b>	<b>R Data, Part 1: Vectors</b>	<b>21</b>
2.1	Vectors	21
2.1.1	Creating Vectors	21
2.1.2	Sequences	22
2.1.3	Logical Vectors	23
2.1.4	Vector Operations	24
2.1.5	Names	27
2.2	Data Types	27
2.2.1	Some Less-Common Data Types	28
2.2.2	What Type of Vector Is This?	28
2.2.3	Converting from One Type to Another	29
2.3	Subsets of Vectors	31
2.3.1	Extracting	31
2.3.2	Vectors of Length 0	34
2.3.3	Assigning or Replacing Elements of a Vector	35
2.4	Missing Data (NA) and Other Special Values	36
2.4.1	The Effect of NAs in Expressions	37
2.4.2	Identifying and Removing or Replacing NAs	37
2.4.3	Indexing with NAs	39
2.4.4	NaN and Inf Values	40
2.4.5	NULL Values	40
2.5	The <code>table()</code> Function	40
2.5.1	Two- and Higher-Way Tables	42
2.5.2	Operating on Elements of a Table	42
2.6	Other Actions on Vectors	45
2.6.1	Rounding	45
2.6.2	Sorting and Ordering	45
2.6.3	Vectors as Sets	46
2.6.4	Identifying Duplicates and Matching	47
2.6.5	Finding Runs of Duplicate Values	49
2.7	Long Vectors and Big Data	50
2.8	Chapter Summary and Critical Data Handling Tools	50
<b>3</b>	<b>R Data, Part 2: More Complicated Structures</b>	<b>53</b>
3.1	Introduction	53
3.2	Matrices	53
3.2.1	Extracting and Assigning	54
3.2.2	Row and Column Names	56

3.2.3	Applying a Function to Rows or Columns	57
3.2.4	Missing Values in Matrices	59
3.2.5	Using a Matrix Subscript	60
3.2.6	Sparse Matrices	61
3.2.7	Three- and Higher-Way Arrays	62
3.3	Lists	62
3.3.1	Extracting and Assigning	64
3.3.2	Lists in Practice	65
3.4	Data Frames	67
3.4.1	Missing Values in Data Frames	69
3.4.2	Extracting and Assigning in Data Frames	69
3.4.3	Extracting Things That Aren't There	72
3.5	Operating on Lists and Data Frames	74
3.5.1	Split, Apply, Combine	75
3.5.2	All-Numeric Data Frames	77
3.5.3	Convenience Functions	78
3.5.4	Re-Ordering, De-Duplicating, and Sampling from Data Frames	79
3.6	Date and Time Objects	80
3.6.1	Formatting Dates	80
3.6.2	Common Operations on Date Objects	82
3.6.3	Differences between Dates	83
3.6.4	Dates and Times	83
3.6.5	Creating POSIXt Objects	85
3.6.6	Mathematical Functions for Date and Times	86
3.6.7	Missing Values in Dates	88
3.6.8	Using Apply Functions with Dates and Times	89
3.7	Other Actions on Data Frames	90
3.7.1	Combining by Rows or Columns	90
3.7.2	Merging Data Frames	91
3.7.3	Comparing Two Data Frames	94
3.7.4	Viewing and Editing Data Frames Interactively	94
3.8	Handling Big Data	94
3.9	Chapter Summary and Critical Data Handling Tools	96
<b>4</b>	<b>R Data, Part 3: Text and Factors</b>	<b>99</b>
4.1	Character Data	100
4.1.1	The <code>length()</code> and <code>nchar()</code> Functions	100
4.1.2	Tab, New-Line, Quote, and Backslash Characters	100
4.1.3	The Empty String	101
4.1.4	Substrings	102
4.1.5	Changing Case and Other Substitutions	103
4.2	Converting Numbers into Text	103
4.2.1	Formatting Numbers	103

4.2.2	Scientific Notation	106
4.2.3	Discretizing a Numeric Variable	107
4.3	Constructing Character Strings: Paste in Action	109
4.3.1	Constructing Column Names	109
4.3.2	Tabulating Dates by Year and Month or Quarter Labels	111
4.3.3	Constructing Unique Keys	112
4.3.4	Constructing File and Path Names	112
4.4	Regular Expressions	112
4.4.1	Types of Regular Expressions	113
4.4.2	Tools for Regular Expressions in R	113
4.4.3	Special Characters in Regular Expressions	114
4.4.4	Examples	114
4.4.5	The <code>regexr()</code> Function and Its Variants	121
4.4.6	Using Regular Expressions in Replacement	123
4.4.7	Splitting Strings at Regular Expressions	124
4.4.8	Regular Expressions versus Wildcard Matching	125
4.4.9	Common Data Cleaning Tasks Using Regular Expressions	126
4.4.10	Documenting and Debugging Regular Expressions	127
4.5	UTF-8 and Other Non-ASCII Characters	128
4.5.1	Extended ASCII for Latin Alphabets	128
4.5.2	Non-Latin Alphabets	129
4.5.3	Character and String Encoding in R	130
4.6	Factors	131
4.6.1	What Is a Factor?	131
4.6.2	Factor Levels	132
4.6.3	Converting and Combining Factors	134
4.6.4	Missing Values in Factors	136
4.6.5	Factors in Data Frames	137
4.7	R Object Names and Commands as Text	137
4.7.1	R Object Names as Text	137
4.7.2	R Commands as Text	138
4.8	Chapter Summary and Critical Data Handling Tools	140
<b>5</b>	<b>Writing Functions and Scripts</b>	<b>143</b>
5.1	Functions	143
5.1.1	Function Arguments	144
5.1.2	Global versus Local Variables	148
5.1.3	Return Values	149
5.1.4	Creating and Editing Functions	151
5.2	Scripts and Shell Scripts	153
5.2.1	Line-by-Line Parsing	155
5.3	Error Handling and Debugging	156
5.3.1	Debugging Functions	156

5.3.2	Issuing Error and Warning Messages	158
5.3.3	Catching and Processing Errors	159
5.4	Interacting with the Operating System	161
5.4.1	File and Directory Handling	162
5.4.2	Environment Variables	162
5.5	Speeding Things Up	163
5.5.1	Profiling	163
5.5.2	Vectorizing Functions	164
5.5.3	Other Techniques to Speed Things Up	165
5.6	Chapter Summary and Critical Data Handling Tools	167
5.6.1	Programming Style	168
5.6.2	Common Bugs	169
5.6.3	Objects, Classes, and Methods	170
<b>6</b>	<b>Getting Data into and out of R</b>	<b>171</b>
6.1	Reading Tabular ASCII Data into Data Frames	171
6.1.1	Files with Delimiters	172
6.1.2	Column Classes	173
6.1.3	Common Pitfalls in Reading Tables	175
6.1.4	An Example of When <code>read.table()</code> Fails	177
6.1.5	Other Uses of the <code>scan()</code> Function	181
6.1.6	Writing Delimited Files	182
6.1.7	Reading and Writing Fixed-Width Files	183
6.1.8	A Note on End-of-Line Characters	183
6.2	Reading Large, Non-Tabular, or Non-ASCII Data	184
6.2.1	Opening and Closing Files	184
6.2.2	Reading and Writing Lines	185
6.2.3	Reading and Writing UTF-8 and Other Encodings	187
6.2.4	The Null Character	187
6.2.5	Binary Data	188
6.2.6	Reading Problem Files in Action	190
6.3	Reading Data From Relational Databases	192
6.3.1	Connecting to the Database Server	193
6.3.2	Introduction to SQL	194
6.4	Handling Large Numbers of Input Files	197
6.5	Other Formats	200
6.5.1	Using the Clipboard	200
6.5.2	Reading Data from Spreadsheets	201
6.5.3	Reading Data from the Web	203
6.5.4	Reading Data from Other Statistical Packages	208
6.6	Reading and Writing R Data Directly	209
6.7	Chapter Summary and Critical Data Handling Tools	210

- 7 Data Handling in Practice 213**
  - 7.1 Acquiring and Reading Data 213
  - 7.2 Cleaning Data 214
  - 7.3 Combining Data 216
    - 7.3.1 Combining by Row 216
    - 7.3.2 Combining by Column 218
    - 7.3.3 Merging by Key 218
  - 7.4 Transactional Data 219
    - 7.4.1 Example of Transactional Data 219
    - 7.4.2 Combining Tabular and Transactional Data 221
  - 7.5 Preparing Data 225
  - 7.6 Documentation and Reproducibility 226
  - 7.7 The Role of Judgment 228
  - 7.8 Data Cleaning in Action 230
    - 7.8.1 Reading and Cleaning `BedBath1.csv` 231
    - 7.8.2 Reading and Cleaning `BedBath2.csv` 236
    - 7.8.3 Combining the `BedBath` Data Frames 238
    - 7.8.4 Reading and Cleaning `EnergyUsage.csv` 239
    - 7.8.5 Merging the `BedBath` and `EnergyUsage` Data Frames 242
  - 7.9 Chapter Summary and Critical Data Handling Tools 245
  
- 8 Extended Exercise 247**
  - 8.1 Introduction to the Problem 247
    - 8.1.1 The Goal 248
    - 8.1.2 Modeling Considerations 249
    - 8.1.3 Examples of Things to Check 249
  - 8.2 The Data 250
  - 8.3 Five Important Fields 252
  - 8.4 Loan and Application Portfolios 252
    - 8.4.1 Layout of the `Beachside Lenders` Data 253
    - 8.4.2 Layout of the `Wilson and Sons` Data 254
    - 8.4.3 Combining the Two Portfolios 254
  - 8.5 Scores 256
    - 8.5.1 Scores Layout 256
  - 8.6 Co-borrower Scores 257
    - 8.6.1 Co-borrower Score Examples 258
  - 8.7 Updated `KScores` 259
    - 8.7.1 Updated `KScores` Layout 259
  - 8.8 Loans to Be Excluded 260
    - 8.8.1 Sample Exclusion File 260
  - 8.9 Response Variable 260
  - 8.10 Assembling the Final Data Sets 262

8.10.1	Final Data Layout	262
8.10.2	Concluding Remarks	263
<b>A</b>	<b>Hints and Pseudocode</b>	<b>265</b>
A.1	Loan Portfolios	265
A.1.1	Things to Check	266
A.2	Scores Database	267
A.2.1	Things to Check	268
A.3	Co-borrower Scores	269
A.3.1	Things to Check	270
A.4	Updated KScores	271
A.4.1	Things to Check	272
A.5	Excluder Files	272
A.5.1	Things to Check	272
A.6	Payment Matrix	273
A.6.1	Things to Check	274
A.7	Starting the Modeling Process	275
	<b>Bibliography</b>	<b>277</b>
	<b>Index</b>	<b>279</b>



## About the Authors

**Samuel E. Buttrey** received a bachelor's degree in statistics from Princeton University in 1983. After 8 years as a Wall Street computer systems analyst, he returned to graduate school and received MA and PhD degrees in statistics from the University of California at Berkeley, the latter in 1996. In that year, he joined the faculty of the Department of Operations Research at the Naval Postgraduate School in Monterey, California. He has published papers on nearest-neighbor and other classification methods and on applied problems ranging from numismatics and oceanography to human vision. He has also published papers describing his implementations of algorithms in software. His interests include classification, computationally intensive methods, and statistical graphics, and most recently, inter-point distance measures for mixed categorical and numeric data. He lives in Pacific Grove, California, with wife Elinda, son John, and some cats.

**Lyn R. Whitaker** received a bachelor's degree in genetics in 1978 and a PhD in statistics from the University of California, Davis, in 1985. She was an Assistant Professor in the Department of Statistics and Applied Probability at the University of California at Santa Barbara from 1985 to 1988, and joined the faculty of the Department of Operations Research at the Naval Postgraduate School in 1988. Her interests are applied statistics relevant to defense issues. These include unsupervised methods for large and messy data, the statistical aspects of reliability and survival analysis, and most recently, jointly with Buttrey, development and use of inter-point distances for mixed data types. She resides in Monterey, California, with husband Mike, father Fred, and, occasionally, children Alex, Lee, and Mary.



## Preface

Statisticians use data to build models, and they use models to describe the world and to make predictions about what will happen next. There has been a large number of very good books that describe statistical modeling, but these modeling efforts usually start with a set of “clean,” well-behaved data in which nothing is missing or anomalous.

In real life, data is messy. There will be missing values, impossible values, and typographical errors. Data is gathered from multiple sources, leading to both duplication and inconsistency. Data that should be categorical is coded as numeric; data that should be numeric can appear categorical; data can be hidden inside free-form text; and data can be in the form of dates in a wide number of possible formats. We estimate that 80% of the time taken in any data analysis problem is taken up just in reading and preparing the data. So, any analyst needs to know how to acquire data and how to prepare it for modeling, and the steps taken should be automatic, as far as possible, and reproducible.

This book describes how to handle data using the R software. R is the most widely used software in statistics, and it has the advantage of being free, open-source, and available on every major computing platform. Whatever software you use, you will find yourself facing the issues of acquiring, cleaning, and merging data, and documenting the steps you took. We hope this book will help you do these things efficiently.

Monterey, California, USA  
November 30, 2016

Sam Buttrey and Lyn Whitaker



## Acknowledgments

Our book is about how to use R to process data. We use R because it is powerful, versatile, and extensible. We thank the developers of R for their service to the statistical community for producing a high-quality open-source piece of software. We also thank the long list of colleagues and students who have helped frame our thinking about questions of statistics and data.



## About the Companion Website

Don't forget to visit the companion website for this book:

[www.wiley.com/go/buttrey/datascientistsguide](http://www.wiley.com/go/buttrey/datascientistsguide)



There you will find valuable material designed to enhance your learning, including:

- A complete listing of all the R code in the Book
- Example datasets used in the Exercises



## 1

## R

## 1.1 Introduction

This book focuses on one problem that is common to almost every statistical problem – indeed, to almost any problem involving any sort of analysis. That problem is acquiring and preparing the data. Across our many years of data analysis, we have learned that seemingly 80% of our time – maybe more – goes into the data preparation steps (a belief echoed by others such as Dasu and Johnson, 2003). Collectively, we call these actions *data cleaning*, although, as we will discuss later, we sometimes use that term for something a little more specific. Regardless of the name, almost any analysis requires that you (i) acquire that data, that is, read it into the computer program; (ii) clean the data, that is, identify entries that are duplicated or clearly erroneous or anomalous, and take other preparation steps (e.g., combining entries such as “Female,” “female,” and “F”); (iii) merge data from different sources; and (iv) prepare the data for modeling, which might involve dividing a set of numeric values into subsets, combining states into regions, and so on. This book discusses some approaches for accomplishing these four steps in the R language (R Core Team, 2013). A fifth problem, which receives less emphasis, is the problem of long-term curation of the data. Which parts of the data must be saved and in what way? We address that question by reference to the idea of *reproducible research*, which we discuss later in this chapter, and later in the book as well.

### 1.1.1 What Is R?

R is a computer program that lets you analyze data. By “analyze” we mean, first, read the data into the program and then operate on it – drawing graphs and charts, manipulating values, fitting statistical models, and so on. (Notice that we prefer to call data “it” rather than “them.” We discuss this choice briefly toward the end of the chapter.) R is both a statistical “environment” and also

a programming language, and it is very widely used both in commercial and academic settings. R is free and open-source and runs on Windows, Apple, and Linux operating systems. It is maintained by a group of volunteers who release bug fixes and new features regularly.

### 1.1.2 Who Uses R and Why?

R started as a tool for statisticians, evolving from a language called S that was created in the 1970s. Today, R remains the primary language of academic statisticians, and it also has a prominent place among analysts in business and government as well. It is used not only for building statistical models but also for handling and cleaning data, as in this book, and for developing new statistical methods, building simulations, for visualization, and generally for all the data-handling tools the statistician and the data scientist require. Because of the ease with which users can develop and distribute new methods, R has also become the tool of choice in certain fast-growing fields such as biostatistics and genetics. Articles on “surveys of the top tools used by data scientists” inevitably name R as one of the important tools with which data scientists, as well as statisticians, should be familiar. Moreover, R’s popularity is such that there are extensions to R (see “packages” in Section 1.4.4) that allow you to connect to other programs such as the Python and Java languages, the H2O machine-learning system, the ArcGIS geographical information system, and many more.

### 1.1.3 Acquiring and Installing R

The primary way to acquire R is to download it from the Internet. The main R website for R is [www.r-project.org](http://www.r-project.org), and the [www.cran.r-project.org](http://www.cran.r-project.org) page (“CRAN” standing for “Comprehensive R Archive Network”) is where you can download R itself. There are in fact dozens of “mirror” sites for CRAN – that is, websites that are essentially copies of the CRAN site – so as to reduce the load on the CRAN site. You can probably find a mirror near you on the “mirrors” page. After you download R, install it in the way you would normally install a program on your operating system.

At any one time, users around the world will be running slightly different versions of R, since new ones are released fairly frequently. For example, at this writing the current version of R was called 3.3.2, but many users are still using 3.2 or earlier versions. This will almost never cause problems, but it is a good idea to update your version of R from time to time.

There are also several slightly different versions of R distributed other than at CRAN. Microsoft R Open is a particular version of R that uses a different set of math libraries intended to make certain computations faster. Like “regular” R, Microsoft R Open is free, although it does not run on OS X. Other versions of R are intended to communicate with relational databases or with other

big-data platforms. For this book, we will assume you are running “regular” R – but in any case for our purposes all versions of R should behave exactly the same way.

### 1.1.4 Starting and Quitting R

The way you start R depends on your operating system. Normally double-clicking on an R icon will be enough to get R started. In the command-line interface of many Linux systems, or using the OS X terminal window, it may be enough just to type the upper-case letter R (or, for Windows command lines, `Rgui`). When R has started, you will see the command prompt `>`. This is the R *console*, the place where commands are entered. At this point, you can start typing commands to R. When it comes time to quit R, you can either “kill” the window in the usual way (for OS X, the red dot, the lightswitch in the top right, or via the File dialog; for Windows, the red X or File dialog) or you can type the `q()` command. In either case, R will then ask you if you want to “Save workspace image.” If you answer “yes” to this question, R will save to the disk any changes you made during the current session, whereas if you answer “no,” R will return its workspace to the condition it was in when R was last started. We almost always want to answer “yes” to this question!

## 1.2 Data

Data is information about the elements of whatever problem we are investigating. Data comes in many forms, but for our purposes it will always be presented in a set of computer-ready values. For example, a database concerning birds might include text about the habits of the birds, numbers giving lengths and weights of the individuals, maps showing migration patterns, images showing the birds themselves, sound recordings of the birds’ calls, and so on. Although they look very different, all of these different pieces of information can be represented in the computer in digital form in one way or another. In this example, one of our primary tasks might be to ensure that each bird’s description is correctly matched with the correct map, image, and song file. Our data analysis projects rarely include data quite so disparate, but in almost every case we need to acquire data, clean it (a process we start to describe in what follows and continue throughout the book), and prepare it for modeling, and in almost every case we expect our data to consist of both numeric and textual values.

### 1.2.1 Acquiring Data

The first step in a data analysis project, of course, is to get the data into R where it can be manipulated. We are old enough to remember the days when this involved typing all the data from the back of a book or journal paper into a

statistics package by hand, but happily this is not necessary today. On the other hand, data now comes in a variety of formats, few of which were created with the convenience of the data scientist in mind. In Chapter 6, we describe some of these common formats and how to use R to read data effectively.

### 1.2.2 Cleaning Data

We “clean” data when we detect (and, in many cases, remove) anomalies. Anomalies will very often be missing values, but they might also be absurd ones, as when people’s ages are reported as 999 or  $-1$ . Sometimes, as in our earlier example, we might have genders reported as “Female,” “female,” and “F” and we want to combine these three values. In the cleaning process we might learn, for example, that one data source produced no data at all in August 2016; this sort of fact will need to be brought to the attention of the data provider. The data cleaning process also involves merging data from different sources, extracting subsets or reshaping the data in some way. All in all, data cleaning is the process of turning raw data, received from one or more providers, into a data set that can be used in visualization, modeling, and decision-making.

In practice these steps are iterative. Our cleaning process not only informs the modeling, but it sometimes leads us to re-acquire the data in a different, more usable form. Similarly, insights from modeling will often lead us to prepare the data in a new and more revealing way – because it is when we model that we often discover anomalies or other interesting attributes of the data.

### 1.2.3 The Goal of Data Cleaning

What a “clean” data set should look like depends on what your goals are. One useful perspective is given by Wickham (2014), who describes what he calls “tidy” data. A tidy data set is rectangular (or tabular); each row describes one unit of analysis (an observation), and each column gives one measurement (a variable). For example, in a data set giving measurements about people, each row would concern itself with a person, and the columns might give height, weight, age, blood type, and so on.

In some problems, it is not immediately clear what the unit of analysis is. For example, imagine data that describes the locations of boats over the course of a month, as recorded by GPS. For some purposes, a “tidy” data set would have one row per GPS ping, each row giving a ship identifier, a location, and a time. For other purposes, we might prefer a data set with one row per boat, each row giving the southernmost point that ship reaches, or perhaps giving a binary indicator of whether the ship did, or did not, spend time in international waters. Some data – images and sound, for example – do not lend themselves to this “tidy” approach.

The exact layout of your final data will depend on what you plan to do with it – and in some cases this won't be known until after you have operated on the data.

### 1.2.4 Making Your Work Reproducible

It is vital that other people be able to reproduce the actions you took on your data. Ideally, you or another analyst should be able to start with your raw data, run all the steps you applied to it, and emerge with exactly the same clean, prepared data sets. This will be useful to you when you encounter a situation similar to the one in the previous paragraph, where the form of the new data needs to be designed. But it is even more important for another analyst, since if you or another analyst can reproduce your results there will be no disagreement about the data. The act of making research reproducible has, in recent years, been rightfully recognized as a cornerstone of scientific progress. Record and document every step you take so that others can repeat them.

## 1.3 The Very Basics of R

This book is about handling data in R. It cannot teach you the very basics of R in detail – although, happily, there are many good books and online resources that can. (We give a few examples at the end of this chapter.) In this section, we list a few of the most basic facts about R, but, again, this book is not intended to teach you R. Rather, we focus on the details of R and of the way data is represented in R, in order to help you understand some of the ways to acquire, clean, and handle data inside R.

### 1.3.1 Top Ten Quick Facts You Need to Know about R

In this section, we give a few of the most important facts about R a beginner needs to know. There will be more detail on these facts later in the chapter and throughout the book.

- 1) The **prompt** is (by default) `>`. If you leave a command incomplete, maybe because there is an unclosed parenthesis or quotation mark, R gives you the continuation prompt, which is `+`. The Esc key (Windows) or control-C (other systems) produces the **break** command, which will take you back to the regular prompt. In this example, we show what a completed command looks like – in this case, R is computing the value of 3 divided by 2.

```
> 3 / 2
[1] 1.5
```

Here, R produced the prompt (`>`), and we typed `3 / 2` and pressed the Enter (or “Return”) key. R then produced the output. We will talk about the `[1]` part in Chapter 2, but the computed value of `1.5` is shown. In the following example, we show what happens when we press Enter after typing the slash character:

```
> 3/
+ 2
[1] 1.5
```

Here, since the expression on the first line was incomplete, R produced the continuation prompt, `+`. When we typed `2` and hit Enter, the expression was complete and the result was shown. In case of confusion, press break until the original `>` prompt is showing.

In examples in this book where we want to show the R output, we also show the `>` prompt in front of our code. Remember, that `>` is produced by R; you don't need to type that yourself. (At the end of the chapter, we tell you where you can get all the code from the book in electronic form.)

- 2) R is **case-sensitive**, which means that upper- and lower-case letters are different in R. For example, the built-in R object `LETTERS` gives all 26 upper-case letters. A different item called `letters` contains the lower-case versions of the alphabet. There is no built-in object called `Letters`.
- 3) **Show** an object by typing its name. For example, if you type `ls` by itself, you see the contents of the function whose name is `ls`, the one that lists all the objects in your workspace (which we define later). To actually run the function and see the objects, you need to type the function's name together with parentheses. In this case, list your objects by typing `ls()`.
- 4) Get **help** for a function or object named `thing` with the command `help(thing)` or `?thing`. For example, to see the help for the `ls()` function, type `help(ls)`. If you don't know the name, try `help.search()` with a relevant word in quotation marks; for example, try `help.search("matrices")` to see functions that handle matrices.
- 5) **Assign** a value or object to a name with the left-arrow (less-than plus hyphen): for example, the command `a <- 1` creates a new object named `a` with value `1`. (You can also assign with a command such as `a = 1`, but we don't recommend it.) The assignment will over-write any existing object named `a` you might have had. Once you create an object, it is in your “workspace,” and your workspace can be saved when you quit. So unless your computer crashes, when you create an object it will persist until you delete it. **Display** the set of objects in your workspace with `objects()` or `ls()`; remove an object with `remove()` or `rm()`. Not every character is permitted in the name of an R object. Start a name