# **CONTENTS**

| р | rn | fa | ~~ |  |
|---|----|----|----|--|
| М | re | īa | ce |  |

| Pr | eface |   | xi |
|----|-------|---|----|
| 1  | Intro | duction   | 1  |
|    | 1.1   | Book Overview                                     | 3  |
|    | 1.2   | Chapter Summaries                                 | 4  |
|    | 1.3   | How to Use This Book                              | 5  |
|    | 1.4   | Why Learn to Analyze Data?                        | 6  |
|    |       | 1.4.1 Learning to Code                            | 6  |
|    | 1.5   | Getting Ready                                     | 7  |
|    | 1.6   | Introduction to R                                 | 8  |
|    |       | 1.6.1 Doing Calculations in R                     | 9  |
|    |       | 1.6.2 Creating Objects in R                       | 10 |
|    |       | 1.6.3 Using Functions in R                        | 12 |
|    | 1.7   | Loading and Making Sense of Data                  | 14 |
|    |       | 1.7.1 Setting the Working Directory               | 15 |
|    |       | 1.7.2 Loading the Dataset                         | 15 |
|    |       | 1.7.3 Understanding the Data                      | 16 |
|    |       | 1.7.4 Identifying the Types of Variables Included | 19 |
|    |       | 1.7.5 Identifying the Number of Observations      | 20 |
|    | 1.8   | Computing and Interpreting Means                  | 21 |
|    |       | 1.8.1 Accessing Variables inside Dataframes       | 21 |
|    |       | 1.8.2 Means                                       | 22 |
|    | 1.9   | Summary   | 24 |
|    | 1.10  | Cheatsheets                                       | 25 |
|    |       | 1.10.1 Concepts and Notation                      | 25 |
|    |       | 1.10.2 R Symbols and Operators                    | 26 |
|    |       | 1.10.3 R Functions                                | 26 |
| 2  | Estir | nating Causal Effects with Randomized Experiments | 27 |
|    | 2.1   | Project STAR                                      | 27 |
|    | 2.2   | Treatment and Outcome Variables                   | 28 |
|    |       | 2.2.1 Treatment Variables                         | 29 |
|    |       | 2.2.2 Outcome Variables                           | 29 |
|    | 2.3   | Individual Causal Effects                         | 29 |
|    | 2.4   | Average Causal Effects                            | 33 |
|    |       | 2.4.1 Randomized Experiments and the              |    |
|    |       | Difference-in-Means Estimator                     | 35 |
|    | 2.5   | Do Small Classes Improve Student Performance? .   | 39 |

|   |       | 2.5.1 Relational Operators in R                     | 39  |
|---|-------|---|-----|
|   |       | 2.5.2 Creating New Variables                        | 40  |
|   |       | 2.5.3 Subsetting Variables                          | 42  |
|   | 2.6   | Summary   | 46  |
|   | 2.7   | Cheatsheets   | 47  |
|   |       | 2.7.1 Concepts and Notation                         | 47  |
|   |       | 2.7.2 R Symbols and Operators                       | 50  |
|   |       | 2.7.3 R Functions                                   | 50  |
| 3 | Infor | ring Population Characteristics via Survey Research | 51  |
| 5 | 3.1   | The FU Referendum in the UK                         | 51  |
|   | 3.1   | Survey Research                                     | 52  |
|   | 5.2   | 321 Random Sampling                                 | 53  |
|   |       | 322 Potential Challenges                            | 54  |
|   | 33    | Measuring Support for Brevit                        | 55  |
|   | J.J   | 331 Predicting the Referendum Outcome               | 56  |
|   |       | 3.3.2 Froquency Tables                              | 57  |
|   |       | 3.3.3 Tables of Propertiens                         | 57  |
|   | 34    | Who Supported Brovit?                               | 58  |
|   | 5.1   | 341 Handling Missing Data                           | 50  |
|   |       | 342 Two-Way Frequency Tables                        | 62  |
|   |       | 343 Two-Way Tables of Proportions                   | 64  |
|   |       | 344 Histograms                                      | 66  |
|   |       | 345 Densitu Histograms                              | 68  |
|   |       | 346 Descriptive Statistics                          | 71  |
|   | 35    | Relationship between Education and the Leave        | , , |
|   | 5.5   | Vote in the Entire UK                               | 76  |
|   |       | 3.5.1 Scatter Plots                                 | 78  |
|   |       | 3.5.2 Correlation                                   | 82  |
|   | 3.6   | Summaru   | 88  |
|   | 3.7   | Cheatsheets   | 90  |
|   |       | 3.7.1 Concepts and Notation                         | 90  |
|   |       | 3.7.2 R Sumbols and Operators                       | 96  |
|   |       | 3.7.3 R Functions                                   | 96  |
|   | P     |   |     |
| 4 | Pred  | licting Outcomes Using Linear Regression            | 98  |
|   | 4.1   | GDP and Night-Time Light Emissions                  | 98  |
|   | 4.2   | Predictors, Observed vs. Predicted Outcomes, and    | 00  |
|   | 4.2   |   | 99  |
|   | 4.3   | Summarizing the Relationship between Two            | 100 |
|   |       | Variables with a Line                               | 100 |
|   |       | 4.3.1 The Linear Regression Widdel                  | 101 |
|   |       | 4.3.2 The Intercept Coefficient                     | 103 |
|   |       | 4.5.5 The Stope Coemclent                           | 104 |
|   | 1 1   | 4.5.4 The Least Squares Method                      | 100 |
|   | 4.4   | A 4.1 Palationahin batu was CDP and D 1 CDP         | 107 |
|   |       | 4.4.1 Kelationship between GDP and Prior GDP        | 109 |
|   | 4 ⊑   | 4.4.2 With Natural Logarithm Transformations .      | 113 |
|   | 4.5   | Fredicting GDP Growth Using Night-Time Light        | 110 |
|   |       | Emissions   | 110 |

|   | 4.6  | Measu    | ring How Well the Model Fits the Data with   |     |
|---|------|----------|--|-----|
|   |      | the Co   | efficient of Determination, $R^2$            | 120 |
|   |      | 4.6.1    | How Well Do the Three Predictive Models      |     |
|   |      |          | in This Chapter Fit the Data?                | 122 |
|   | 4.7  | Summa    | ary  | 123 |
|   | 4.8  | Append   | dix: Interpretation of the Slope in the Log- |     |
|   |      | Log Lir  | near Model                                   | 124 |
|   | 4.9  | Cheats   | heets  | 126 |
|   |      | 4.9.1    | Concepts and Notation                        | 126 |
|   |      | 4.9.2    | R Functions                                  | 128 |
| 5 | Fati | matina ( | Causal Effects with Observational Data 1     | 120 |
| 5 | 51   | Russia   | n State-Controlled TV Coverage of 2014       | 25  |
|   | 5.1  | Ukrain   | ian Affairs                                  | 129 |
|   | 52   | Challer  | nges of Estimating Causal Effects with       | 125 |
|   | 5.2  | Ohserv   | rational Data                                | 130 |
|   |      | 521      | Confounding Variables                        | 130 |
|   |      | 522      | Why Are Confounders a Problem?               | 131 |
|   |      | 523      | Confounders in Randomized Experiments        | 133 |
|   | 53   | The Ff   | fect of Russian TV on Ukrainians' Voting     | 100 |
|   | 0.0  | Behavi   | or   | 135 |
|   |      | 5.3.1    | Using the Simple Linear Model to Compute     |     |
|   |      |          | the Difference-in-Means Estimator            | 136 |
|   |      | 5.3.2    | Controlling for Confounders Using a          |     |
|   |      |          | Multiple Linear Regression Model             | 142 |
|   | 5.4  | The Ef   | fect of Russian TV on Ukrainian Electoral    |     |
|   |      | Outcon   | nes  | 147 |
|   |      | 5.4.1    | Using the Simple Linear Model to Compute     |     |
|   |      |          | the Difference-in-Means Estimator            | 149 |
|   |      | 5.4.2    | Controlling for Confounders Using a          |     |
|   |      |          | Multiple Linear Regression Model             | 151 |
|   | 5.5  | Interna  | l and External Validity                      | 153 |
|   |      | 5.5.1    | Randomized Experiments vs.                   |     |
|   |      |          | Observational Studies                        | 153 |
|   |      | 5.5.2    | The Role of Randomization                    | 154 |
|   |      | 5.5.3    | How Good Are the Two Causal Analyses         |     |
|   |      |          | in This Chapter?                             | 155 |
|   |      | 5.5.4    | How Good Was the Causal Analysis in          |     |
|   |      |          | Chapter 2?                                   | 156 |
|   |      | 5.5.5    | The Coefficient of Determination, $R^2$      | 157 |
|   | 5.6  | Summa    | ary  | 157 |
|   | 5.7  | Cheats   | heets  | 159 |
|   |      | 5.7.1    | Concepts and Notation                        | 159 |
|   |      | 5.7.2    | R Functions                                  | 161 |
| 6 | Dral | abilitu  | 1  | 162 |
| U | 61   | What I   | s Probabilitu?                               | 162 |
|   | 6.2  | Avione   | of Probabilitu                               | 162 |
|   | 6.3  | Fvente   | Random Variables and Probability             | 105 |
|   | 0.5  | Distrib  | utions                                       | 165 |
|   |      | Distrib  |  | .05 |

|     | 6.4               | Probability Distributions  | 166<br>166<br>169<br>173               |
|-----|-------------------|--|--|
|     | 6.5               | <ul> <li>6.4.4 Recap</li> <li>Population Parameters vs. Sample Statistics</li> <li>6.5.1 The Law of Large Numbers</li> <li>6.5.2 The Central Limit Theorem</li> <li>6.5.3 Sampling Distribution of the Sample Management of t</li></ul> | 179<br>179<br>180<br>183               |
|     | 6.6<br>6.7<br>6.8 | 0.5.3       Sampling Distribution of the Sample Mean         Summary   | 189<br>190<br>192<br>192<br>194<br>195 |
| 7   | Quai              | ntifying Uncertainty   | 196                                    |
|     | 7.1               | Estimators and Their Sampling Distributions  | 196                                    |
|     | 7.2               | Confidence Intervals   | 202<br>203<br>206<br>200               |
|     | 7.3               | Hypothesis Testing   | 203<br>211<br>218<br>220               |
|     | 7.4               | Statistical vs. Scientific Significance  | 224                                    |
|     | 7.5<br>7.6        | Summary  | 225<br>226<br>226<br>229               |
|     |                   | 7.6.3 R Functions  | 229                                    |
| Ind | ex of             | Concepts   | 231                                    |
| Ind | ex of             | Mathematical Notation  | 235                                    |
| Ind | ex of             | R and RStudio  | 237                                    |

# **1. INTRODUCTION**

This book provides a friendly introduction to data analysis for the social sciences. It covers the fundamental methods of quantitative social science research, using plain language and assuming absolutely no prior knowledge of the subject matter.

Proceeding step by step, we show how to analyze real-world data using the statistical program R for the purpose of answering a wide range of substantive questions. Along the way, we teach the statistical concepts and programming skills needed to conduct and evaluate social scientific studies. We explain not only how to perform the analyses but also how to interpret the results and identify the analyses' strengths and potential limitations.

Through this book, you will learn how to *measure*, *predict*, and *explain* quantities of interest based on data. These are the three fundamental goals of quantitative social science research. (See outline 1.1.)

#### WHY DO WE ANALYZE DATA IN THE SOCIAL SCIENCES?

In the social sciences we analyze data to:

- *measure* a quantity of interest, such as the proportion of eligible voters in favor of a particular policy
- predict a quantity of interest, such as the likely winner of an upcoming election
- *explain* a quantity of interest, such as the causal effect of attending a private school on student test scores.

Figuring out whether you aim to measure, predict, and/or explain a quantity of interest should always precede the analysis and often also precede the data collection. As you will learn, the goals of your research will determine (i) what data you need to collect and how, (ii) the statistical methods you use, and (iii) what you pay attention to in the analysis. As you read this book and learn about each goal in detail, the distinctions will become clearer. Here we provide a brief preview. R symbols, operators, and functions introduced in this chapter: +, -, \*, /, <-, ", (), sqrt(), #, setwd(), read.csv(), View(), head(), dim(), \$, and mean().

OUTLINE 1.1. The three goals of quantitative social science research.

To measure a quantity of interest such as a population characteristic, we often use survey data, that is, information collected on a sample of individuals from the target population. To analyze the data, we may compute various descriptive statistics, such as mean and median, and create visualizations like histograms and scatter plots. The validity of our conclusions depends on whether the sample is representative of the target population. To measure the proportion of eligible voters in favor of a particular policy, for example, our conclusions will be valid if the sample of voters surveyed is representative of *all* eligible voters.

To predict a quantity of interest, we typically use a statistical model such as a linear regression model to summarize the relationship between the predictors and the outcome variable of interest. The stronger the association between the predictors and the outcome variable, the better the predictive model will usually be. To predict the likely winner of an upcoming election, for example, if economic conditions are strongly associated with the electoral outcomes of candidates from the incumbent party, we may be able to use the current unemployment rate as our predictor.

To explain a quantity of interest such as the causal effect of a treatment on an outcome, we need to find or create a situation in which the group of individuals who received the treatment is comparable, in the aggregate, to the group of individuals who did not. In other words, we need to eliminate or control for all confounding variables, which are variables that affect both (i) the likelihood of receiving the treatment and (ii) the outcome variable. For example, when estimating the causal effect of attending a private school on student test scores, family wealth is a potential confounding variable. Students from wealthier families are more likely to attend a private school and also more likely to receive after-school tutoring, which might have a positive impact on their test scores. To produce valid estimates of causal effects, we may conduct a randomized experiment, which eliminates all confounding variables by assigning the treatment at random. In the current example, we would achieve this by using a lottery to determine which students attend private schools and which do not. Alternatively, if we cannot conduct a randomized experiment and need to rely on observational data instead, we would need to use statistical methods to control for all confounding variables such as family wealth. Otherwise, we would not know what portion of the difference in average test scores between private and public school students was the result of the type of school attended and what portion was the result of family background.

INTRODUCTION 3

## **1.1 BOOK OVERVIEW**

The book consists of seven chapters.

Chapter 1 is the introductory chapter, which lays the groundwork for the forthcoming data analyses.

Chapters 2 through 5 each introduce one or two published social scientific studies. In these chapters, we show how to analyze real-world datasets to answer different kinds of substantive questions. Specifically, we teach how to use several quantitative methods to measure, predict, and explain quantities of interest. (See outline 1.2, which indicates how each chapter relates to the three goals of quantitative social science research.)

|    | BOOK OUTLINE   |           |
|----|--|-----------|
| Ch | apter  | Goal      |
| 1. | Introduction   |           |
| 2. | Estimating Causal Effects with<br>Randomized Experiments | Explain   |
| 3. | Inferring Population Characteristics via Survey Research | Measure   |
| 4. | Predicting Outcomes Using Linear<br>Regression           | Predict   |
| 5. | Estimating Causal Effects with<br>Observational Data     | Explain   |
| 6. | Probability  |           |
| 7. | Quantifying Uncertainty                                  | All Three |

As you can see, chapters 2 and 5 are both about explanation, also known as causal inference. They teach how to estimate causal effects using different types of data. Since the methods differ, they are presented in separate chapters.

The book progresses from simple to more complex methods. Chapter 2 shows how to estimate causal effects using data from a randomized experiment. Chapter 3 is about measurement and teaches how to infer the characteristics of an entire population from a sample of survey respondents. Chapter 4 is about prediction and demonstrates how to use simple linear regression. Chapter 5 shows how to estimate causal effects with observational data and teaches multiple linear regression, the most complicated method we see in the book.

In chapter 6, we cover basic probability, and in chapter 7 we complete some of the analyses from chapters 2 through 5 by quantifying the uncertainty of our empirical findings. A more detailed description of each chapter is below.

OUTLINE 1.2. Book outline showing how each chapter relates to the three goals of quantitative social science research.

## **1.2 CHAPTER SUMMARIES**

In the current introductory chapter, we discuss why data analysis is a required skill among social scientists. We also explain how to get our computers ready, and we familiarize ourselves with RStudio and R, the two programs we will use. Then, we learn to load and make sense of data and practice computing and interpreting means.

In chapter 2, we define and learn how to estimate causal effects using data from a randomized experiment. As the working example, we analyze data from one of the largest experiments in U.S. education policy research, Project STAR, to determine whether attending a small class improves student performance.

In chapter 3, we use survey research to measure population characteristics. In addition, we learn how to visualize and summarize the distribution of single variables as well as the relationship between two variables. To illustrate these concepts, we analyze data related to the 2016 British referendum on withdrawing from the European Union, a decision popularly known as Brexit.

In chapter 4, we learn how to predict outcomes using simple linear regression models. For practice, we analyze data from 170 countries in order to predict growth in gross domestic product (GDP) using night-time light emissions as measured from space.

In chapter 5, we return to estimating causal effects, but this time using observational data. We define confounding variables, examine how their presence complicates the estimation of causal effects, and learn how to use multiple linear regression models to help mitigate the potential bias these variables introduce. To illustrate how this works step by step, we estimate the effects of Russian TV reception on the 2014 Ukrainian parliamentary elections. In this context, we introduce the concepts of internal and external validity. We then discuss the pros and cons of randomized experiments and of observational studies.

In chapter 6, we shift our focus away from data analysis to cover basic probability. We learn about random variables and their distributions as well as the distinction between population parameters and sample statistics. We then discuss the two large sample theorems that enable us to measure statistical uncertainty.

In chapter 7, we use everything we have learned in the preceding chapters and show how to quantify the uncertainty in our empirical findings in order to draw conclusions at the population level. In particular, we show how to quantify the uncertainty in (i) population inferences, (ii) predictions, and (iii) causal effect estimates. As illustrations, we complete some of the analyses we started in chapters 2 through 5.

## 1.3 HOW TO USE THIS BOOK

This is no ordinary textbook on data analysis. It is intentionally designed to accommodate readers with a variety of math and programming backgrounds.

The book uses a two-column layout: a main column and a side column or margin.

The main column contains the essential material and code, which are intended for all readers, except for the sections labeled FORMULA IN DETAIL. These contain more advanced material and are clearly identified so that you can easily skip them if you so choose.

In the margin are various types of notes and figures, each with a different purpose:

- At the beginning of each chapter, we list the R functions, symbols, and operators that will be introduced. You can look through the list to get a sense of what will be covered. (See, for example, the list for this chapter shown on the first page, and note that we always display code in cyan.)
- TIPs include supplemental material, such as additional explanations, answers to common questions, notes on best practices, and recommendations.
- RECALLs remind you of relevant information mentioned earlier in the book. These reminders are particularly helpful when the book is read only a few pages at a time, such as over the course of a semester.
- To help you review the core concepts, which are shown in **bold red** in the main text, we repeat their definitions in the margin. These notes are displayed in **red**.
- To help you with R functions, symbols, and operators, the first time these are introduced, we include in the margin an explanation of how they work and provide an example. These explanations are displayed in a cyan-colored frame.

At the end of each chapter, in place of the usual list of supplementary exercises, we include CHEATSHEETS to help you review the core concepts as well as the R functions, symbols, and operators covered.

Supplementary chapter-specific exercises, categorized by degree of difficulty, are available at http://press.princeton.edu/dss.

Finally, at the end of the book, we include three separate indexes for concepts, mathematical notation, and R-related topics.

## 1.4 WHY LEARN TO ANALYZE DATA?

As a social scientist, sooner or later you will need to rely on data to (i) measure the characteristics of a certain population of interest, (ii) make predictions, and/or (iii) make or evaluate decisions involving cause-and-effect relationships. What proportion of a population is in favor of a particular policy? Who is the candidate most likely to win an upcoming election? Shall we implement a particular policy to boost economic growth? You will want to be able to answer these types of questions either by analyzing data yourself or by understanding and assessing someone else's data analysis.

Even if you are not planning to become a social scientist, it is useful for you to know how to analyze data and/or how to distinguish a good quantitative study from a poorly conducted one. These are highly marketable skills. Recent advancements in computing power and the proliferation of data have increased the demand for data analysts who can inform decision makers in the public and private sectors alike.

The analytical skills you will learn by making your way through this book can also be used to improve everyday decisions, from choosing a candidate to vote for to determining the best way to increase your productivity. Perhaps most importantly, by learning the strengths and limitations of different quantitative methods, you will become less vulnerable to arguments based on faulty inferences from data. In the era of big data, we all stand to benefit from becoming savvy consumers of quantitative research, even if we do not all become skilled researchers ourselves.

#### 1.4.1 LEARNING TO CODE

For the purpose of analyzing data, we write and run code. Code contains instructions that a computer can implement. These instructions consist of sequences of clearly defined steps written in a particular programming language. In this book, we code in R, which is a programming language used by many data analysts.

Don't worry if you have never done any coding before. Learning to code is not as difficult as one might think. You may even find it fun. Back in 1944, when the first programmable computer in the United States was built, only highly trained mathematicians were able to code. At that time, coding required punching paper tape in specific sequences that the machine could read. (See a rendition of what this tape looked like in the margin.) Today, anyone with access to a computer, some spare time, and a little patience can learn how to code.



INTRODUCTION 7

## 1.5 GETTING READY

To perform the analyses in this book, we first need to download and install the necessary files and programs. We should also familiarize ourselves with RStudio, which is the interface we use throughout.

#### **1** DOWNLOAD AND SAVE FILES

All the files we will use are in a folder named DSS, which is available at <u>http://press.princeton.edu/dss</u>. For easy access, we recommend saving the folder on your Desktop. This is where the code used throughout the book assumes the DSS folder is located. In case you choose to save the folder elsewhere, we also provide instructions for making the necessary changes to the code.

#### **2** DOWNLOAD AND INSTALL R AND RSTUDIO

We will use two programs: R and RStudio. R is the statistical program, the engine if you will, that will perform the calculations and create the graphics for us. RStudio is the user-friendly interface we will use to communicate with R. While we could use R directly, going through RStudio makes writing and running code much easier.

Why do we use R as our statistical program? Because it is free, open-source (anyone can see the underlying code and improve it), powerful, and flexible. It is also widely used. Indeed, many jobs these days require knowledge of R.

TIP: By default, your computer will likely save the DSS folder to your Downloads. To move it, you can copy and paste it or drag it to the new location.

Unfortunately, these programs are compatible only with Linux, Mac, and Windows operating systems. They cannot be used on tablets or phones. We provide instructions for using these two programs on a Mac or a Windows computer.



To download and install R, go to <u>http://cran.r-project.org</u>, select the link that matches your operating system, and follow the instructions.



To download and install RStudio, go to <u>http://rstudio.com</u>, select the link that matches your operating system, and follow the instructions.

#### **③** BECOME FAMILIAR WITH RSTUDIO

To analyze data, we always operate R through RStudio. Let's take a moment to become acquainted with RStudio's layout.

After installing both programs, go ahead and start RStudio. Then, from within RStudio, open a new R script, which is the type of file we use to store the code we write to analyze data. Instructions are shown in the margin.

TIP: How do we open a new R script? In the RStudio dropdown menu, click on File > New File > R Script. A new "Untitled" file will open. The extension of this type of file is ".R", which is why R scripts are also called R files. After opening a new R script, RStudio's interface should look like figure 1.1.

- The upper-left window is the *R script*, which is where we write and run code, giving R commands to execute.
- The lower-left window is the *R* console, where R provides either the results of successfully executed code (known as outputs) or any error messages.
- The upper-right window is the *environment*, which is the storage room of the current R session. It lists all the objects we have created. (We will soon explain what objects are and provide examples showing how the environment works.)
- The lower-right window is where we find the *help* and *plots tabs*, which we will learn how to use later on.

| RStudio File Edit Code View Plots Session Bu | ld Debug Profile Tools Window Help       |
|--|--|
| ⊡ Untitled ×                                 | Environment History                      |
| R SCRIPT                                     | ENVIRONMENT                              |
| Console   Terminal ×                         | Files   Plots   Packages   Help   Viewer |
| R CONSOLE                                    |  |

## **1.6 INTRODUCTION TO R**

To use R, we need to learn the R programming language. (R is the name of both the statistical program and the programming language.) Learning a programming language is like learning a foreign language. It is not easy, and it takes a lot of practice and patience. The exercises in this book will help you learn to code in R, so be sure to follow along. Practice is everything!

Let's begin. R can be used to do many things. In our case, we will use R (i) as a calculator; (ii) to create objects, which is how R stores data; and (iii) to interact with data using functions.

#### WE WILL USE THE STATISTICAL PROGRAM R TO:

- (i) do calculations
- (ii) create objects
- (iii) use functions.

FIGURE 1.1. Layout of RStudio after opening a new R script. The upper-left window is the R script. The lower-left window is the R console. The upper-right window is the environment of the R session. The plots and help tabs appear in the lower-right window.

INTRODUCTION 9

#### 1.6.1 DOING CALCULATIONS IN R

We can use R as a calculator. R can do summation (+), subtraction (-), multiplication  $(^*)$ , and division (/), as well as other more complicated mathematical operations. For example, the code to ask R to calculate 1 plus 3 is:

#### 1 + 3

To run this or any other code, we first type it in the R script (the upper-left window of RStudio). Then, we highlight as much of it as we want to run and either (a) manually hit the run icon (shown in the margin) or (b) use the shortcut *command+enter* in Mac or *ctrl+enter* in Windows. The result, or output, of the executed code will show up in the R console (the lower-left window of RStudio). (Instead, we could type the code directly in the R console and hit enter, but we should avoid doing it that way. It is best to run code through an R script so that you can save it, re-run it, tweak it, expand it, and share it.)

After running the code above, we should see the following in the R console: first, the executed code shown in blue, indicating that R was able to run it without problems, and then the output shown in black. In this case, the output is:

4

Indeed, one plus three equals four.

Congratulations! You just wrote and ran your first line of code in R. Notice that now that you have written some code in the R script, RStudio shows the name of the file in red. This is to remind you that you have some unsaved changes. Once you save the file, the file name will return to black.

Throughout the book, we show the output that you should see in the R console right after the code that produces it. To distinguish the output from the code, we display the output with the symbol ## at the beginning of the line. For example, we display the code and output above as follows:

1 + 3 ## [1] 4

The first line, shown in cyan, is the code to be typed and run in the R script. The second line, which begins with ## and is shown in gray, is what should appear in the R console after running the code.

What does the number in brackets before the 4 mean? It indicates the position of the output immediately to its right. In this instance, [1] indicates that 4 is the first output of the code we ran. Later in the chapter, we will see examples of code that produce multiple outputs, which will clarify how this works. +, -, \*, and / are some of the arithmetic operators recognized by R. Example: (4 - 1 + 3) \* (2 / 3)



TIP: To save any changes you make to the R script, either (a) use the shortcut *command*+S in Mac or *ctrl*+S in Windows or (b) click on File > Save or Save As...

TIP: Adding spaces around operators makes the code easier to read. R ignores these spaces. Example: 1+3 produces the same output as 1 + 3.



<- is the assignment operator. It creates new objects in R (unless one with the same name already exists, in which case R overwrites its contents). To its left, we specify the name of the object (without quotes). To its right, we specify the contents of the object. Example: four <- 4.

### 1.6.2 CREATING OBJECTS IN R

In order to manipulate and analyze data, we need to load and store datasets. R stores information in what are known as objects, and so we need to learn how to create objects in R.

Think of an object as a box that can contain anything. All we need to do is give it a name, so that we know how to refer to it, and specify its contents.

To create an object in R, we use the assignment operator <-:

- To its left, we specify the name we want to give the object. This name can be anything as long as it does not begin with a number or contain spaces or special symbols like \$ or % that are reserved for other purposes. Underscores \_ are permitted and are good substitutes for spaces.
- To its right, we specify the contents of the object, that is, the data we want to store.

CREATING OBJECTS: To store data as an object in R, we run code using this format:

#### object\_name <- object\_contents</pre>

where:

- object\_name is the name we want to give the object
- <- is the assignment operator, which creates an object by assigning contents to a name
- *object\_contents* is the data we want to store in the object.

For example, if we want to create an object called *four* containing the output of the calculation 1+3, we run:

four <-1+3

Notice that after running the code above, the object will show up in the environment (the upper-right window in RStudio). As mentioned earlier, the environment is the storage room of the current R session. It shows the objects that we have created and that are available for us to use.

If we want to know the contents of the object *four*, we can type and run the name of the object in the R script. Its contents will appear in the R console. This is equivalent to asking R, what is inside the object named *four*?

**four** ## [1] 4

TIP: We would accomplish the same thing by running: four <-4.

TIP: RStudio continues to work in the same R session until you quit the program. At that time, R will ask whether you want to save the workspace image, which contains all the objects you created during the R session. We recommend that you do not save it. If you need to continue to work with those objects, you can always re-create them by re-running your code.

Not surprisingly, the object four contains the number 4.

Objects can contain text as well as numbers. For example, to create an object called *hello* containing the text "hi" we run: hello <- "hi"

After running the code above, the environment should contain two objects: *four* and *hello*.

Let's stop here to learn something important about R. Look at the code above. Why did we use quotation marks around the content of the object "hi" but not around the name of the object hello? In other words, when do we use quotes " when coding in R? Here is the rule: When writing code, the names of objects, names of functions, and names of arguments as well as special values such as TRUE, FALSE, NA, and NULL should *not* be in quotes; all other text should be in quotes. (In the next subsection, we will see what we mean by functions and arguments. We will learn the meaning and usage of TRUE and FALSE in chapter 2 and of NA and NULL in chapter 3.)

What would have happened had we tried to run the code above without quotes around *hi*? Go ahead and try it:

hello <- hi
## Error: object 'hi' not found</pre>

In the R console, you will see an error message (in red) that reads, "Error: object 'hi' not found". Indeed, by typing *hi* without quotes, you are telling R that *hi* is the name of an object. Because there is no object named *hi* in the environment, R gives you an error message. Encountering programming errors is part of the coding process. Try not to be discouraged by them.

A word of caution: R overwrites (replaces) old objects if we use the same name when creating a new object. For example, go ahead and run the following:

hello <- "hi, nice to meet you"

You should see that you still have only two objects in the environment: *four* and *hello*, but now *hello* contains the text "hi, nice to meet you" instead of simply "hi". To confirm this, we run:

hello
## [1] "hi, nice to meet you"

Note also that R is case-sensitive. It will treat *Hello* as a completely different object name than *hello*. If we run the name *Hello* by mistake, R will not be able to find the object because there is no object in the environment called *Hello* with an uppercase H at the beginning. To avoid this problem, we recommend using all lowercase letters when naming objects. " when writing code, the names of objects, names of functions, and names of arguments as well as special values such as TRUE, FALSE, NA, and NULL should not be in quotes; all other text should be in quotes. Examples: "this is just text", object\_name. Never use quotes around a number unless you want R to treat it as text, in which case you will not be able to use it to perform arithmetic operations.

TIP: If you have problems figuring out what a particular error means, Google it. Lots of data analysts participate in Q&A sites, such as Stack Overflow, which can be very helpful for this sort of thing.

#### 1.6.3 USING FUNCTIONS IN R

Finally, we use R to interact with data, which requires using functions.

Think of a function as an action that you request R to perform with a particular piece of data, such as calculating the square root of four. A function takes one or more inputs, such as the number four; performs one or more actions with the inputs, such as calculating the square root of the inputs; and produces one or more outputs, such as the number two, which is the result of taking the square root of four.



R functions in this book: sqrt(), setwd(), read.csv(), View(), head(), dim(), mean(), ifelse(), table(), prop.table(), na.omit(), hist(), median(), sd(), var(), plot(), abline(), cor(), lm(), log(), c(), sample(), rnorm(), pnorm(), print(), nrow(), predict(), abs(), and summary().

() the names of functions are always followed by parentheses. Inside the parentheses, we write the argument(s) of the function, separated by commas if there is more than one argument. Example: function\_name(arg1, arg2).

TIP: There are two types of arguments: required and optional. Required arguments are the inputs that we must specify in order to use a particular function. Optional arguments are the inputs that we could specify if we wanted to modify the function's default settings. Throughout the book, we will learn how to use the functions listed in the margin, which come automatically loaded with R. In time, we will learn their names, the actions they perform, the inputs they require, and the outputs they produce. Meanwhile, here are some important things to know about functions:

- The name of a function (without quotes) is always followed by parentheses: *function\_name()*
- Inside the parentheses, we specify the inputs to be used by the function, which we refer to as arguments: function\_name(arguments)
- Most functions require that we specify at least one argument but can take many optional arguments. When multiple arguments are specified inside the parentheses, they are separated by commas: *function\_name(argument1, argument2)*
- To identify the type of argument that we are specifying, we either enter the arguments in a particular order or include their names (without quotes) in our specification: function\_name(argument1, argument2) or function\_name(argument1\_name = argument1, argument2\_name = argument2)
- In this book, we follow the most common practices. We always specify required arguments first. If there is more than one required argument, we enter them in the order expected by R. We specify any optional arguments we want next and include their names so that R knows how to interpret them: function\_name(required\_argument,

optional\_argument\_name = optional\_argument)

INTRODUCTION 13

USING R FUNCTIONS: To use a function in R, we typically write code in one of these two formats:

(a) function\_name(required\_argument)

(b) function\_name(required\_argument, optional\_argument\_name = optional\_argument)

where:

- *function\_name* is the name of the function; for example, "mean" is the name of the function that computes the mean of a set of values
- required\_argument is the argument the function requires, such as the values we want to calculate the mean of; we typically do not include the names of required arguments; we enter the required arguments first, and if there is more than one, we enter them in the order expected by R
- , is a comma, which we use to separate different arguments
- optional\_argument\_name is the name of the optional argument we want to use, such as the argument that enables us to eliminate missing values before calculating a mean
- *optional\_argument* is what we set the optional argument to be.

We will see some complex R functions in the next section. For now, let's look at a simple one. The function sqrt(), which stands for "square root," calculates the square root of the argument specified inside the parentheses. For example, to calculate the square root of 4, we run:

**sqrt (4)** ## [1] 2

The output here is the number 2. Alternatively, given that the object *four* currently contains the number 4, we can run:

**sqrt (four)** ## [1] 2

Note that R will be able to execute the code above only after we have created the object *four*. If we start a new R session and attempt to run this code without having first created the object *four*, R will not be able to find the object in the environment and will give us an error message. This is just to say that one must run code in order. When returning to work on an R script, it is a good idea to run all the code from the beginning of the file up to the line that we are working on. sqrt() calculates the square root of the argument specified inside the parentheses. Example: sqrt(4).

TIP: Here, the name of the function is sqrt, which, as with all function names, is followed by parentheses (). Inside the parentheses, we need to specify the required argument, which is 4, in this case. The output of the executed code is 2. Indeed, the square root of four is two.

One of the major advantages of writing code using an R script (instead of writing it directly into the R console) is that we can always replicate our results by re-running the code we have written previously. Using an R script, we are able to work on complex problems that might require running hundreds or thousands of lines of code. As long as we save the code in an R script, we can keep tweaking and expanding it. Writing code in R scripts also means we can share our work and collaborate with others. Anyone with access to our R script will be able to replicate our analyses, which leads us to our next topic: the importance of annotating, or commenting, code.

It is good practice to comment code, that is, to include short notes to ourselves or to our collaborators explaining what the code does. This will make reading and understanding our code easier. To write comments in the R script, we use #. R ignores everything that follows this character until the end of the line and will not execute it. For example, running the following code produces exactly the same output as the code above:

```
sqrt(four) # calculates square root of four
## [1] 2
```

After seeing the *#* character, R stopped reading until the end of the line. Had we inserted the comment at the beginning of the line, before the code, R would not have executed the function sqrt() at all. Go ahead, run:

# calculates square root of four sqrt(four)

R will not produce an output. R thinks the whole line is a comment because it starts with #.

RStudio helps us write and read code by color coding it in the R script. For example, comments are shown in light green, while executable code is shown in black, gray, blue, and dark green. Becoming familiar with this color scheme will help you detect errors in your code. (In this book, we use only two colors when displaying code: cyan for executable code and gray for comments.)

### 1.7 LOADING AND MAKING SENSE OF DATA

Before starting any analysis, we must load the dataset. Then, we must understand what the observations represent and what each of the variables means. In this section, we show how to do all of this for the data from the Project Student-Teacher Achievement Ratio (Project STAR) in preparation for the analysis in chapter 2. The goal of Project STAR was to examine the causal effects of small classes on student performance. While exploring the

# is the character used to comment code. R ignores everything that follows this character until the end of the line. Example: # this is a comment.

TIP: We recommend that when you start a new study, you either (a) start a new R session (Session > New Session) or (b) remove all objects from the environment (Session > Clear Workspace > Yes) to avoid operating with objects from previous studies by mistake. data from Project STAR, we learn what variables are and how to distinguish between different types of variables based on their contents.

To follow along, you can create a new R script (as shown in the previous section) and practice typing the code yourself. Alternatively, you can open the "Introduction.R" file, which contains the code used in the remainder of this chapter.

Here are the steps we recommend you follow before starting a data analysis:

#### SET THE WORKING DIRECTORY

Before we can load a dataset, we need to direct R to the working directory, that is, the name and location of the folder containing the data. If you followed the advice from the earlier section, all the files necessary for the exercises in this book will be in the DSS folder on your Desktop.

The easiest way to set the working directory is to first save the R script to the folder that contains the dataset, the DSS folder, in this case. Then, you can use the dropdown menu to set the working directory manually: Session > Set Working Directory > To Source File Location. After your last click, you will see a line of code appear in the R console. Every time you start a new R session and want to work with a dataset saved in the DSS folder, you will need to run this line of code. You may, therefore, want to copy and paste it in the R script as your first line of code.

The code to set the working directory uses the function setwd(), which stands for "<u>set working directory</u>." The only required argument is the path to the folder, which should be in quotes because it is text and not the name of an object, the name of a function, the name of an argument, or a special value such as TRUE, FALSE, NA, and NULL. The path differs depending on whether you have a Mac or a Windows computer. The code to set the working directory to the DSS folder on your Desktop should resemble one of these (where *user* is your own username):

setwd("~/Desktop/DSS") # example of setwd() for Mac
setwd("C:/user/Desktop/DSS") # example for Windows

#### **2** LOAD THE DATASET

Now we are ready to load the dataset. R can read a variety of data formats. In this book, datasets are always provided in comma-separated values files, known as CSV files. As the name indicates, CSV files contain data separated by commas. (See figure 1.2 for a rendition of a CSV file.) TIP: How do we open an existing R script? In the RStudio dropdown menu, click on File > Open File... and then click on the ".R" file you want to open.

TIP: To save an R script to the DSS folder, either (a) click on File > Save As ... and select the DSS folder, or (b) manually drag the corresponding ".R" file from its current location to the DSS folder.

setwd() sets the working directory, that is, directs R to the folder on your computer where the dataset is saved. The only required argument is the path to the folder in quotes. Examples: setwd("~/Desktop/folder") for Mac, setwd("C:/user/Desktop/folder") for Windows (where user is your own username).

TIP: Resist the temptation to double-click on a CSV file. If you open a CSV file directly, you risk inadvertently changing or losing data.

16 CHAPTER 1

FIGURE 1.2. CSV files contain data separated by commas.

read.csv() reads CSV files. The only required argument is the name of the CSV file in quotes. Example: read.csv("file.csv").

| ○ ○ ○ □ STAR.csv                            |
|---|
| "classtype", "reading", "math", "graduated" |
| "small",578,610,1                           |
| "regular",612,612,1                         |
| "regular",583,606,1                         |
| "small",661,648,1                           |
| "small",614,636,1                           |
| "regular",610,603,0                         |
|   |

To read the contents of a CSV file in R, we use the read.csv() function, which requires that we specify inside the parentheses the name of the CSV file in quotes. (We need to use quotes around the name of the CSV file because it is text and not the name of an object, the name of a function, the name of an argument, or a special value such as TRUE, FALSE, NA, and NULL.)

To store the dataset so that we can analyze it later, we need to not only read the CSV file but also save its contents as an object. We can do so by using the assignment operator <-. As we saw earlier, to the left of the assignment operator, we specify the name of the object. To its right, we specify the contents, which in this case are produced by reading the CSV file using the function read.csv().

Here, the dataset is in a file called "STAR.csv", and we choose to name the object where we store the dataset *star*. Putting it all together, the code to read and store the dataset is:

star <- read.csv("STAR.csv") # reads and stores data</pre>

After running the line of code above, the name of the object, *star*, should appear in the environment (the upper-right window in RStudio). If R gives you an error message instead, make sure that (i) you have set the working directory to the folder where the CSV file is saved, (ii) the name of the CSV file you are using in the code is exactly the same as the name of the CSV file saved in the working directory, and (iii) the extension of the CSV file in the working directory is indeed ".csv".

#### **O UNDERSTAND THE DATA**

To make sense of the dataset, we should start by looking at its contents.

To look at the data, we could type the name of the object, *star*, in the R script and run it. R will show the entire contents of the dataset in the R console, which might be hard to read unless the dataset is small.

A better option is to use the function View(), which requires that we specify inside the parentheses the name of the object where the dataset is stored (without quotes). Alternatively, you can manually click on the object name in the environment. Both of these actions open a new tab in the upper-left window of RStudio with the dataset in spreadsheet form. We can then easily scroll up, down, left, and right to look at the data in an organized manner. Figure 1.3 shows how the data will be displayed if we run:

View(star) # opens a new tab with contents of dataset

| Introduction.R × I star × |           |    |         |     |      |    |                    |        |
|---------------------------|-----------|----|---------|-----|------|----|--------------------|--------|
| •                         | classtype | 40 | reading | 40  | math | 40 | gr <b>a</b> duated | ☆<br>⊽ |
| 1                         | small     |    | 57      | 8'8 | 63   | 10 |                    | 1      |
| 2                         | regular   |    | 61      | .2  | 6    | 12 |                    | 1      |
| 3                         | regular   |    | 58      | 3   | 60   | )6 |                    | 1      |
| 4                         | small     |    | 66      | 51  | 64   | 18 |                    | 1      |
| 5                         | small     |    | 61      | .4  | 63   | 36 |                    | 1      |
| 6                         | regular   |    | 61      | 0   | 60   | )3 |                    | 0      |

Sometimes it might be enough for us to see only the first few rows of data. For this purpose, we use the function head(), which requires that we specify inside the parentheses the name of the object where the dataset is stored (without quotes):

| head( | ( <mark>star)</mark> # s | hows the | first | six rows  |  |
|-------|--------------------------|----------|-------|-----------|--|
| ##    | classtype                | reading  | math  | graduated |  |
| ## 1  | small                    | 578      | 610   | 1         |  |
| ## 2  | 2 regular                | 612      | 612   | 1         |  |
| ## 3  | 3 regular                | 583      | 606   | 1         |  |
| ## 4  | a small                  | 661      | 648   | 1         |  |
| ## 5  | 5 small                  | 614      | 636   | 1         |  |
| ## 6  | o regular                | 610      | 603   | 0         |  |

By default, the function head() displays the first six lines of data. If we want R to display a different number of lines, we specify an optional argument inside the parentheses. In particular, we specify the argument named n and set it to equal the number of lines we want R to show. For example, to ask R to display the first three lines of *star*, we run:

| head(st | tar, n=3 | <mark>3)</mark> | ws the | e first three | rows |
|---------|----------|-----------------|--------|---------------|------|
| ## cl   | asstype  | reading         | math   | graduated     |      |
| ## 1    | small    | 578             | 610    | 1             |      |
| ## 2    | regular  | 612             | 612    | 1             |      |
| ## 3    | regular  | 583             | 606    | 1             |      |

View() opens a new tab in the upper-left window of RStudio with the contents of a dataset. The only required argument is the name of the object where the dataset is stored (without quotes). Example: View(*data*). Note that R is case-sensitive and the name of this function starts with uppercase V. The good news is this is the only function we will see in this book that uses any uppercase letters; all others are written in all lowercase letters.

FIGURE 1.3. Tab that opens in the upper-left window of RStudio with entire contents of the *star* dataset as a result of either (a) running View(star) in the R script or (b) clicking on the object *star* in the environment. (To return to the R script, we can either close this tab by clicking on the gray X next to the name or by clicking on the R script tab.)

head() shows the first six rows or observations in a dataset. The only required argument is the name of the object where the dataset is stored (without quotes). To change the number of observations displayed, we use the optional argument n. Examples: head(*data*) shows the first six rows and head(*data*, n=3) shows the first three. In the output, the first column identifies the position of the observations, and the first row identifies the names of the variables.

RECALL: In R functions, multiple arguments are separated by commas, and the names of arguments should not be in quotes. Failing to follow these instructions will prevent R from executing the code and result in an error message. How do we make sense of the data inside the dataset? Knowing the following common features of datasets should help:

- Datasets capture the characteristics of a particular set of individuals or entities: citizens, organizations, countries, and so on.
   As we will soon learn, this dataset contains information about students who participated in Project STAR.
- Datasets are typically organized as **dataframes**, where rows are observations and columns are variables.
  - What is an observation?
    - An observation is the information collected from a particular individual or entity in the study.
    - The unit of observation of the dataset defines the individuals or the entities that each observation in the dataframe represents. The unit of observation in the STAR dataset is students. Hence, every row of data in the *star* dataframe represents a different student in the study.
    - We usually refer to an observation by the row number in the dataframe, which we denote as *i*. See, for example, that in the output of head(), the rows are labeled by their position, *i*. When R displays the first six observations of a dataframe, the values of *i* range from 1 to 6.
  - What is a variable?
    - A variable captures the values of a changing characteristic for the multiple individuals or entities in the study.
    - Every column of data in the *star* dataframe is a variable; each variable captures a specific feature of the students, for all the students in the study.
    - We usually refer to a variable by its name. See, for example, that in the output of head(), the columns are labeled with the variable names: *classtype, reading, math,* and *graduated.* (Note that, for easy recognition, we italicize the names of variables in the text.)
    - From time to time, in this book we define new variables for the purpose of illustrating concepts. We represent a variable and its contents using the following mathematical notation:

$$X = \{10, 5, 8\}$$

- On the left-hand side of the equal sign, we identify the name of the variable: *X*, in this case.
- On the right-hand side of the equal sign and inside curly brackets, we have the contents of the variable: multiple observations, separated by commas. In the simple example above, the variable contains three observations: 10, 5, and 8.

In a **dataframe**, each row is an observation, and each column is a variable. An **observation** is the information collected from a particular individual or entity in the study. The **unit of observation** of a dataset defines what each observation represents. The **notation** *i* identifies the position of the observation; the observation for which i=1 is the first observation. A **variable** captures the values of a changing characteristic for multiple individuals or entities.

TIP: In this book, we are also going to teach you mathematical notation, a system of symbols and expressions that represent mathematical concepts. To help you keep track of the meaning of these symbols and expressions, at the end of the book we have included an index of all the mathematical notation we use.

- To represent each individual observation of the variable X, we use  $X_i$  (pronounced X sub i) where i stands for the observation number. The subscript i means that we have a different value of X for each value of i. Here, because there are only three observations, i can equal only 1, 2, or 3. For example, we represent the second observation (the observation for which i=2) as  $X_2$ , which in this case equals 5.

Now that we have looked at the data, we should read the description of the variables provided in table 1.1.

| variable  | description  |
|-----------|--|
| classtype | class size the student attended: "small" or<br>"regular"                                     |
| reading   | student's third-grade reading test scores (in points)  |
| math      | student's third-grade math test scores (in points)   |
| graduated | identifies whether the student graduated from high school: 1=graduated or 0=did not graduate |

Reading table 1.1, we learn that:

- classtype captures the size of the class the student attended, which was either "small" or "regular"
- *reading* records the scores, measured in points, the student earned on the third-grade reading test
- math records the scores, measured in points, the student earned on the third-grade math test
- *graduated* indicates whether the student graduated from high school (it equals 1 if the student graduated and 0 if the student did not graduate).

Now we can look at the first few lines of data again and substantively interpret them. For example, the first observation represents a student who attended a small class, earned 578 points on the third-grade reading test and 610 points on the third-grade math test, and graduated from high school.

#### **O** IDENTIFY THE TYPES OF VARIABLES INCLUDED

At this point, we should learn the typology of the variables included in the dataset. This information will be especially helpful when we need to interpret the results of the analysis.

Based on the contents of the variables, we can distinguish between the types listed in outline 1.3.

TABLE 1.1. Description of the variables in the STAR dataset, where the unit of observation is students.

OUTLINE 1.3. Types of variables based on their contents.

A character variable contains text, such as names={Elena, Kosuke, Kathryn}. A numeric variable contains numbers, such as rank= $\{2, 1, 3\}$ . A binary variable can take only two values; in this book, we define binary variables as taking only 1s and 0s, such as voted= $\{1, 0, 1\}$ . A non-binary variable can take more than two values, such as distance= $\{1.452, 2.345, 0.298\}$  and dice\_roll= $\{2, 4, 6\}$ .



dim() provides the dimensions of a dataframe. The only required argument is the name of the object where the dataframe is stored (without quotes). The output is two values: the first indicates the number of observations in the dataframe; the second indicates the number of variables. Example: dim(*data*).



- The first distinction we make is between character and numeric variables. While character variables contain text, numeric variables contain numbers. For example, in the STAR dataset, *classtype* is a character variable, and *reading*, *math*, and *graduated* are numeric variables.
- Among numeric variables, we differentiate between binary and non-binary. A binary variable can take only two values ("bi" means two). In this book, all binary variables take only 1s and Os to represent the presence or absence of a particular trait. In this type of binary variable, also known as a dummy variable, you may think of the 1s as capturing the positive responses to simple yes or no questions and of the 0s as capturing the negative responses. For example, in the STAR dataset, graduated is a binary variable that captures responses to the question, did the student graduate? The variable takes the value of 1 when the answer is yes (the student graduated) and 0 when the answer is no (the student did not graduate). (See mathematical definition in the margin.)
- In contrast, we categorize as non-binary all other numeric variables, that is, those that can take more than two values.
   For example, in the STAR dataset, both *reading* and *math* are non-binary variables because they each contain more than two different numbers.

#### **IDENTIFY THE NUMBER OF OBSERVATIONS**

Finally, we should find out how many observations the dataset contains. For this purpose, we use the function dim(), which stands for "dimensions" and requires that we specify inside the parentheses the name of the object where the dataframe is stored (without quotes). This function returns two values because dataframes have two dimensions: rows and columns. The first corresponds to the number of rows, which is equivalent to the number of observations. The second corresponds to the number of columns, which is equivalent to the number of to the number of variables.

INTRODUCTION 21

Given the output above, we learn that the STAR dataset contains 1,274 observations. And, as we already knew by looking at the data directly, it contains four variables. Given that each observation represents a student, we now know that we have information for 1,274 students in Project STAR. In mathematical notation, we represent the number of observations in a dataframe as n. In this case, we can state that n=1,274.

#### **1.8 COMPUTING AND INTERPRETING MEANS**

The mean, or average, of a variable is one of the foundational concepts of data analysis. In this section, we first show how to access a variable inside a dataframe in R so that we can operate with its values. Then, we explain in detail how to calculate and interpret the mean of a variable.

#### **1.8.1 ACCESSING VARIABLES INSIDE DATAFRAMES**

Suppose we want to operate with the variable *reading* inside the dataframe *star*. How do we access the values within this variable?

If we run the name of the variable, *reading*, R will give us an error message informing us that the object *reading* cannot be found. Indeed, there is no object called *reading* in the environment. If instead we run the name of the object that contains the dataframe, *star*, R will show all the values in the dataframe, not just those of the variable *reading*.

To access the values of a single variable, we use the \$ character. To its left, we specify the name of the object where the dataframe is stored (without quotes). To its right, we specify the name of the variable (without quotes). In this case, *star*\$*reading* is the code that instructs R to select the variable *reading* from within the object named *star*. It is equivalent to saying to R: look inside of *star* and find the variable called *reading*. (Note that when writing code, we do not use quotes around the names of elements within an object, such as the names of variables within a dataframe.) Go ahead and run:

star \$reading

```
## [1] 578 612 583 661 614 610 595 665 616 624
## [11] 593 599 693 545 565 654 686 570 529 582
## ...
```

In your R console (the lower-left window), you should see all the observations of *reading*. Here, we show you only the first 20. We use an ellipsis—three dots—to signify that more observations should appear after those displayed here. The **notation** *n* stands for the total number of observations in a dataframe or in a variable.

TIP: In chapter 3, we will see how to calculate and interpret other statistics, such as the median, standard deviation, and variance of a variable.

\$ is the character used to access an element inside an object, such as a variable inside a dataframe. To its left, we specify the name of the object where the dataframe is stored (without quotes). To its right, we specify the name of the variable (without quotes). Example: *data*\$*variable* accesses the variable named *variable* inside the dataframe stored in the object named *data*.

TIP: The number in brackets shown on the second line in your R console might not be 11 because the size of your lower-left window might be different than ours. The larger the window, the more observations R will be able to display per line, and the higher the number in brackets shown on the second line will be.

The **mean**, or **average**, of a variable equals the sum of the values across all observations divided by the number of observations.

TIP: The mean of a variable is a single number, which does not vary by observation. As a result, the mean of  $X(\overline{X})$  is not subscripted by *i*.

mean() computes the mean of a variable. The only required argument is the code identifying the variable. Example: mean(*data*\$*variable*). What are the numbers in brackets at the beginning of each line? They indicate the position of the observation immediately to the right. The [1] on the first line indicates that the number 578 is the first observation of the variable *reading* (*reading*<sub>1</sub>=578). The [11] on the second line indicates that the number 593 is the 11th observation of *reading* (*reading*<sub>11</sub>=593), and so on.

Now that we know how to access the contents of a variable, we can learn how to compute and interpret the mean of a variable.

#### **1.8.2 MEANS**

The mean, or average, of a variable characterizes its central tendency. It equals the sum of the values across all observations divided by the number of observations. In mathematical notation, the mean of a variable is often represented by the name of the variable with a bar on top, like so:

#### name of the variable

The formula to compute the mean of X is:

$$\overline{X} = \frac{\sum_{i=1}^{n} X_i}{n} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

where:

- $\overline{X}$  (pronounced X-bar) stands for the mean of X
- X<sub>i</sub> (pronounced X sub i) stands for a particular observation of X, where *i* denotes the position of the observation
- *n* is the number of observations in the variable
- the symbol  $\sum$  (the Greek letter Sigma) is the mathematical notation for summation;  $\sum_{i=1}^{n} X_i$  stands for the sum of all  $X_i$  (observations of X) from i=1 to i=n, meaning from the first observation of the variable X to the last one.

For example, if  $X = \{10, 4, 6, 8, 22\}$ , then n=5 because the variable has five observations, and the mean of X is:

$$\overline{X} = \frac{\sum_{i=1}^{n} X_i}{n} = \frac{X_1 + X_2 + X_3 + X_4 + X_5}{5}$$
$$= \frac{10 + 4 + 6 + 8 + 22}{5} = \frac{50}{5} = 10$$

To calculate the mean of a variable in R, we can use the function mean(). The only required argument is the code identifying the variable. For example, to calculate the mean of the reading test scores of the students in the STAR dataset, we run:

mean(star\$reading) # calculates the mean of reading
## [1] 628.803

How shall we interpret this output? First, we need to figure out the quantity in which the value is measured. This is called the unit of measurement. When interpreting numeric results, you should make it clear whether the number is measured in points, percentages, miles, or kilometers, for example.

The unit of measurement of the mean of a variable depends on whether the variable is non-binary or binary. Outline 1.4 summarizes how to interpret the mean of a variable (including units of measurement) based on this distinction. (We exclude from this discussion categorical variables, whose means generally have no straightforward substantive interpretation.)

| interpretation of the mean of a variable |                        |  |  |
|--|------------------------|--|--|
| if variable is non-binary:               | if variable is binary: |  |  |
| as an average, in the same               | as a proportion, in %  |  |  |
| unit of measurement                      | after multiplying      |  |  |
| as the variable                          | the result by 100      |  |  |

When the variable is non-binary, the mean should be interpreted in the same unit of measurement as the values in the variable. For example, in the output above, because *reading* is a nonbinary variable measured in points, the mean of *reading* is also in points. We can, therefore, interpret the output as meaning that the students in Project STAR scored 629 points, on average, on the third-grade reading test.

When the variable is binary, the mean should be interpreted as a percentage, after multiplying the result by 100. Why? Because the mean of a binary variable is equivalent to the proportion of the observations that are 1s (that have the characteristic identified by the variable). Let's go over a simple example to see how this works. Suppose we want to calculate the mean of the first six observations of the binary variable graduated. As we saw earlier in the output of head(), these are {1, 1, 1, 1, 1, 0}. The average of these six observations would be:

$$\frac{\text{sum of observations}}{\text{number of observations}} = \frac{1+1+1+1+1+0}{6} = \frac{5}{6} = 0.8333...$$

Notice that the fraction 5/6 is equivalent to the proportion of the observations that are 1s. (See TIP in the margin.) Now, to convert the result from decimal form (0.83) into a percentage, we multiply it by  $100 (0.83 \times 100 = 83\%)$ .

The **unit of measurement** is the quantity in which a value is measured. For example, depending on where you live, you might measure temperature in  $^{\circ}F$  or  $^{\circ}C$  and distance in miles or kilometers.

TIP: Categorical, or factor, variables take a fixed number of values, where each value represents a qualitative outcome. For example, we can capture adult education levels in a categorical variable where 1=no qualifications, 2=high school diploma, and 3=undergraduate degree.

OUTLINE 1.4. Interpretation of the mean of a variable based on the type of variable.

TIP: It is good practice to round numeric results to meaningful decimal places. This usually means no more than two decimals, but often one or none.

TIP: The proportion of observations that meet a criterion is calculated as:

number of observations that meet criterion total number of observations

To interpret the resulting decimal as a percentage, we multiply it by 100. For example, if  $X = \{1, 1, 1, 1, 1, 0\}$ , the proportion of observations in X that are 1s is 83% (5/6 = 0.8333... and  $0.83 \times 100 = 83\%$ ).

When the variable is binary, the numerator of the mean, the sum of the 1s and 0s, is equivalent to the number of observations that meet the criterion. As a result, the mean is equivalent to the proportion of observations in the variable that are 1s.

Putting it all together, we interpret the average of the first six observations of *graduated* as indicating that about 83% of the first six students in the STAR dataset graduated.

Now, let's compute the mean of all the observations within the binary variable *graduated* (rather than just the first six). We do so by running:

mean(star\$graduated) # calculates the mean of graduated
## [1] 0.8697017

How shall we interpret this output? Since the variable is binary, the output means that about 87% of all the students in Project STAR received a high school diploma ( $0.87 \times 100 = 87\%$ ).

### 1.9 SUMMARY

We began this chapter by providing an overview of the book and its main features. We then argued that knowing how to perform and evaluate data analyses is particularly useful for studying society and human behavior but can be helpful to anyone. These skills are also highly marketable in the current era of big data. We then got our computers ready and became acquainted with R and RStudio, the two programs we use. Finally, in preparation for the next chapter's data analysis, we learned how to load and make sense of data and how to compute and interpret means.

## **1.10 CHEATSHEETS**

### **1.10.1 CONCEPTS AND NOTATION**

| concept/notation                                | description  | example(s)  |
|---|--|---|
| dataframe                                       | structure of data in which each row is an observation and each column is a variable  | variables <i>classtype</i> and <i>reading</i> in dataframe form:  |
|   | variables  | i classtype reading   |
|   |  | 1 small 725   |
|   | observations $\begin{array}{ccc} 1 \rightarrow & & \\ 2 \rightarrow & & \\ \dots & & & \\ \end{array}$   | 2 regular 692<br>3 small 725  |
| observation                                     | information collected from a particular<br>individual or entity in the study; each row<br>in a dataframe is an observation   | the first observation in the dataframe<br>above represents a student who attended<br>a small class and scored 725 points on the<br>reading test |
| unit of observation                             | defines what each observation represents   | students, schools, states, countries,   |
| i   | identifies the position of the observation   | the observation for which $i=3$ is the third observation in a dataframe or in a variable  |
| variable<br>(X)                                 | captures the values of a changing<br>characteristic for multiple individuals or<br>entities; each column in a dataframe is a<br>variable   | <i>X</i> ={10, 8, 12}   |
| character variable                              | variable that contains text  | <i>names</i> = {Elena, Kosuke, Kathryn}   |
| numeric variable                                | variable that contains numbers   | $rank = \{2, 1, 3, 5, 4\}$  |
| binary variable                                 | variable that can take only two values, in<br>this book 1s and Os  | $voted = \{0, 1, 1, 0, 1\}$   |
| non-binary variable                             | variable that can take more than two<br>values   | distance = {2.568, 5.367, 7.235}<br>dice_roll = {2, 5, 4, 3, 1}   |
| п   | stands for the total number of observations in a dataframe or in a variable  | in the variable <i>dice_roll</i> above, <i>n</i> =5   |
| $\sum$  | Greek letter Sigma; mathematical notation  | if X={10, 8, 12}, then:   |
|   | for summation; $\sum_{i=1}^{n}$ means "sum what follows for all the observations, from $i=1$ to $i=n$ " (the first to the last)  | $\sum_{i=1}^{n} X_i = X_1 + X_2 + X_3$<br>= 10 + 8 + 12 = 30  |
| unit of<br>measurement                          | quantity in which a value is measured  | °F, °C, miles, kilometers, points, percentages, percentage points,  |
| mean or average<br>o <u>f</u> a variable<br>(X) | characterizes the central tendency of the<br>variable; equals the sum of the values<br>across all observations divided by the<br>number of observations:                                     | if $X = \{10, 8, 12\}$ , then:<br>$\overline{X} = \frac{\sum_{i=1}^{n} X_i}{n} = (X_1 + X_2 + X_3)/3$   |
|   | mean of $X = \overline{X} = \frac{\sum_{i=1}^{n} X_i}{n}$  | =(10+8+12)/3=30/3=10  |
|   | <ul> <li>unit of measurement of X:</li> <li>if X is non-binary: in the same unit of measurement as X</li> <li>if X is binary: in percentages, after multiplying the result by 100</li> </ul> |   |

### 1.10.2 R SYMBOLS AND OPERATORS

| code    | description  | example(s)   |
|---------|--|--|
| + - * / | some of the arithmetic operators recognized by R   | (4 - 1 + 3) * (2 / 3)  |
| <-      | assignment operator; creates new objects or overwrites existing<br>ones (if an object with the same name already exists); to its left,<br>we specify the name of the object (without quotes); to its right,<br>we specify the contents of the object; the name of an object<br>cannot begin with a number or contain spaces, \$, or %                              | object_name <- object_contents   |
| H       | when writing code, the names of objects, names of functions,<br>and names of arguments as well as special values such as<br>TRUE, FALSE, NA, and NULL should not be in quotes; all<br>other text should be in quotes; numbers should not be in quotes  | "this is text"<br>object_name  |
| ()      | the names of functions are always followed by parentheses;<br>inside the parentheses, we write the argument(s) of the<br>function, separated by commas if there is more than one<br>argument; we enter required arguments first and in the order<br>expected by R; we specify optional arguments by including<br>their names in the specification (without quotes) | function_name(required_argument)<br>function_name(required_argument,<br>optional_argument_name =<br>optional_argument) |
| #       | character used to comment code; R ignores everything that follows this character until the end of the line   | <pre>executable_code # comment</pre>   |
| \$      | character used to access an element inside an object, such as a<br>variable inside a dataframe; to its left, we specify the name of<br>the object where the dataframe is stored (without quotes); to its<br>right, we specify the name of the variable (without quotes)  | <pre>data\$variable # accesses the variable named variable inside the dataframe stored in the object named data</pre>  |

#### 1.10.3 R FUNCTIONS

| function   | description   | required argument(s)  | example(s)  |
|------------|---|---|---|
| sqrt()     | calculates the square root  | what we want to compute the square root of                              | sqrt(4)   |
| setwd()    | sets the working directory,<br>that is, directs R to the<br>folder on your computer<br>where the dataset is saved   | path to folder in quotes  | <pre>setwd("~/Desktop/folder")   # for Mac setwd("C:/user/Desktop/folder")   # for Windows, where user is   your own username</pre> |
| read.csv() | reads CSV files   | name of CSV file containing the dataset in quotes                       | read.csv("file.csv")  |
| View()     | opens a tab with the entire contents of a dataframe   | name of object where the dataframe is stored                            | View( <i>data</i> )   |
| head()     | shows the first six<br>observations in a dataframe;<br>in the output, observations<br>are identified by their<br>position, <i>i</i> , and variables by<br>their names | name of object where the<br>dataframe is stored                         | <pre>head(data) # shows first six observations</pre>  |
|            |   | optional argument n: changes<br>the number of observations<br>displayed | <pre>head(data, n=2) # shows first two observations</pre>   |
| dim()      | provides the dimensions of a<br>dataframe; output is:<br>[1] number of rows or<br>observations<br>[2] number of columns or<br>variables                               | name of object where the<br>dataframe is stored                         | dim( <i>data</i> )  |
| mean()     | calculates the mean of a<br>variable  | code identifying the variable<br>(see \$)                               | mean( <i>data\$variable</i> )   |

# Index of Concepts

absolute value, 219 alternative hypothesis, 212, 218, 222, 227 average, 22, 25 average causal effects, 33, 48 average estimation error, 199, 226 average outcome for the control group, 37, 50 average outcome for the treatment group, 37, 50 average treatment effects, 33, 48 axioms, 163

Bayesian interpretation of probabilities, 162, 192 Bernoulli distribution, 166, 179, 192 bias, 199, 226 binary random variables, 20, 25, 29, 179 binary variables, 166, 192 bivariate linear model, 121

categorical variables, 23 causal assumption, 43, 130, 134, 139, 147, 151, 153 causal effects, 27, 29, 129, 196, 206, 218, 220, 224 causal inference, 27, 129, 211 causal language, 43, 139, 145 causal relationships, 28, 47 central limit theorem, 183, 194 character variables, 20, 25 coefficient of determination, 120, 128, 157 coefficients, 103, 104, 127, 140, 141, 144, 151, 220 confidence interval, 202, 227 confidence level, 202 confounders, 130, 159 confounding variables, 130, 159 control condition, 29, 47 control group, 35, 49 control variables, 159 correlation, 82, 95, 109, 121 correlation coefficient, 82, 95, 121 counterfactual outcomes, 32, 48 critical value, 215, 228 cross-tabulations, 62, 91 cumulative distribution function, 172, 174 dataframes, 18, 25 density function, 170, 193 density histograms, 68, 93 dependent variables, 28, 47, 99, 126 descriptive statistics, 71, 93 difference-in-means estimator, 37, 49, 136, 138, 206, 218 dummy variables, 20 estimate, 196, 226 estimated intercept, 127

estimated slope, 127 estimation error, 199, 226 estimator, 196, 197, 226

232 INDEX OF CONCEPTS

events, 163, 192 expectation, 180, 194 experimental data, 38, 50 experiments, 35, 49 explaining a quantity of interest, 27, 129, 196, 206, 218, 220, 224 external validity, 153, 160 extrapolation, 111

factor variables, 23 factual outcomes, 32, 48 fitted line, 105, 126 fitted linear model, 102, 105, 126, 144 fitted log-log linear model, 116, 128 fitted multiple linear regression model, 160 fitted multiple linear regression model where  $X_1$ is the treatment variable, 160 fitted simple linear regression model, 126 fitted simple linear regression model, 126 fitted simple linear regression model where X is the treatment variable, 159 frequency tables, 57, 90 frequentist interpretation of probabilities, 162, 192

histograms, 66, 92 hypothesis testing, 211, 227 independent variables, 28, 47, 99, 126 individual causal effects, 29, 48 intercept, 103 intercept coefficient, 103, 127 internal validity, 153, 160

large sample theorems, 180 law of large numbers, 180, 194 least squares, 106 level of test, 214, 228 line of best fit, 82, 100, 101, 106 linear association, 82, 95 linear model, 101, 121, 126, 143

item nonresponse, 54, 90

linear regression, 101, 126 linear relationship, 82, 83, 95, 109 log-log linear model, 116, 124 logarithmic transformation, 113

margin of error, 205, 227 mean, 22, 25, 72 measuring a quantity of interest, 51, 196, 203 median, 72, 93 misreporting, 54, 90 multiple linear regression model, 121, 143, 159 mutually exclusive events, 163, 192

natural logarithm, 113 non-binary random variables, 20, 25, 29, 179 non-binary variables, 169, 193 nonlinear relationship, 88 nonresponse, 54, 90 normal distributions, 169, 179, 193 normal random variables, 169, 179 null hypothesis, 211, 218, 222, 227 numeric variables, 20, 25

observational data, 38, 50 observational studies, 38, 50, 153 observations, 18, 25 observed outcomes, 100, 126 omitted variables, 159 one-sided p-value, 214 outcome, 47, 163, 192 outcome variables, 28, 47, 99, 126

p-value, 213, 228
p.p., 44, 49, 58, 127, 138, 140, 147, 150, 151, 153, 205, 222, 227
parameter, 179, 196, 226
percentage change, 44, 117
percentage points, 44, 49, 58, 127, 138, 140, 147, 150, 151, 153, 205, 222, 227
percentage-point change, 44, 117
population mean, 180, 194
population variance, 180

post-treatment variables, 146, 160 potential outcome, 30 potential outcome under the control condition, 47 potential outcome under the treatment condition, 47 pre-treatment characteristics, 36, 49 predicted outcomes, 100, 126, 209 predicting a guantity of interest, 98, 196, 209 prediction errors, 100, 107, 126 predictors, 99, 126 probability, 162 probability density function, 170, 193 probability distributions, 165, 192 proof by contradiction, 211 proportions, 23 random sampling, 53, 90, 154 random treatment assignment, 35, 36, 49, 134, 154 random variables, 165, 192 randomized controlled trials, 35, 49 randomized experiments, 35, 49, 133, 153 RCT, 35, 49 regression line, 101, 126 replication, 203, 215 representative sample, 52, 90 residuals, 100, 107 sample, 52, 90 sample mean, 180, 188, 194, 203 sample space, 163, 192 sample statistic, 179 sample variance, 180, 194 sampling distribution, 188, 194, 198, 226 sampling frame, 54, 90 sampling variability, 180, 194 sampling with replacement, 168, 195 sampling without replacement, 168, 195 scatter plots, 78, 94 scientific significance, 224, 228 significance level, 214, 228 simple linear model, 121

simple linear regression model, 101, 121, 126 simulations, 167, 171, 186, 187 slope, 104 slope coefficient, 104, 127 standard deviation, 73, 93 standard error, 198, 226 standard normal distribution, 173, 193 standardize, 176 standardized estimator, 201 standardized variable, 176, 201 statistical controls, 143 statistical significance, 215, 228 sum of squared residuals (SSR), 107, 120 survey research, 51 surveys, 51 t-distribution, 223 t-statistic, 223 table of proportions, 57, 91 test statistic, 212, 218, 222, 227 total sum of squares (TSS), 120 treatment condition, 29, 47 treatment group, 35, 49 treatment variables, 28, 47 trial, 163, 192 two-sided p-value, 213, 228 two-way frequency tables, 62, 91 two-way tables of proportions, 64, 92 type I error, 215 unbiased estimator, 199, 226 uniform distribution, 67 unit nonresponse, 54, 90 unit of measurement, 23, 25, 43, 112, 138, 140, 141, 151 unit of observation, 18, 25 variables, 18, 25 variance, 76, 94, 194 z-scores, 84, 94, 176 z-statistic, 212, 218, 222, 228

# Index of Mathematical Notation

α, 101, 143, 159  $\hat{\alpha}$ , 103, 127, 140, 144, 160 β, 101  $\beta_i$ , 143, 159  $\widehat{eta}$ , 102, 104, 127, 138, 140, 151, 153 β<sub>1</sub>, 145–147, 220, 222  $\beta_i$ , 144, 160 △, 30, 47 *△X*, 104  $\triangle Y_i$ , 30, 34, 48  $\triangle \hat{Y}$ , 104, 112, 127  $\epsilon_i$ , 101, 143, 159 *ϵ*, 100, 126  $\hat{\epsilon}_{i}$ , 102  $\theta$ , 212, 228 μ, 169–171, 173, 176, 179, 193  $\pi$ , 170, 193  $\sum$ , 22, 25 σ, 169, 170, 173, 176, 193  $\sigma^2$ , 169–171, 173, 176, 179, 193 Ω, 163, 164, 192 ≈, 167 ~, 169, 171, 173, 176, 179, 193  $\sim^{\text{approx.}}$ , 183 →, 28, 47, 130, 142, 145, 159 cor(X, Y), 82, 95, 121

E, 180, 184, 199, 200, 218 **E**(*X*), 180, 181, 183, 188, 194  $\mathbb{E}[Y_i(X_i=1) - Y_i(X_i=0)], 218$ e, 170, 193 estimate<sub>i</sub>, 199, 200 estimation error<sub>i</sub>, 199 estimator, 200-202, 212 H<sub>0</sub>, 211, 227 H<sub>1</sub>, 212, 227 *i*, 18, 25 loq(X), 113, 128 median(X), 72, 93 N (number of observations in the population), 52 N (a normal distribution), 169, 171, 173, 176, 179, 183, 193 name of the variable, 22, 25 N(0, 1), 172, 173, 176, 183, 193, 201, 202, 212 n, 21, 22, 25, 34, 52, 181, 183, 204 n control group, 207, 218 n treatment group, 207, 218 P, 163, 192 p, 166, 179, 192 R<sup>2</sup>, 120, 128, 157

236 INDEX OF MATHEMATICAL NOTATION

sd(X), 73, 93 SSR, 107, 120 (standard error)<sup>2</sup>, 200 standard error of  $\hat{\beta}_1$ , 222 standard error, 200–202, 212

true value, 199–201 *TSS*, 120

 $\mathbb{V}$ , 180, 185, 200  $\mathbb{V}(X)$ , 180, 183, 188, 194 var(X), 76, 94, 180, 194 var(Y), 204  $var(Y_{control})$ , 207, 218  $var(Y_{treatment})$ , 207, 218

X (a predictor), 99, 126 X (a random variable), 18, 22, 25, 166, 169, 171, 176, 179, 181, 183, 188, 193 X (the treatment variable), 28, 47, 130, 159 X<sub>1</sub> (the treatment variable), 145 X<sub>2</sub>,..., X<sub>p</sub> (statistical controls), 145 X<sub>i</sub>, 18, 22, 25, 101, 102

*X<sub>ij</sub>*, 143, 144, 159, 160 <del>X</del>, 22, 25, 72, 180, 181, 183, 188, 194 x, 166, 170 Y, 28, 47, 99, 100, 126, 130, 159 Y<sub>i</sub>, 101, 126, 143, 159  $Y_i(X_i=0)$ , 30, 47, 218  $Y_i(X_i=1)$ , 30, 47, 218  $\overline{Y(X=0)}$ , 36, 48  $\overline{Y(X=1)}$ , 36, 48 <del>Y</del>, 204 Ycontrol group, 37, 49, 50, 207, 218 *Y*<sub>treatment group</sub>, 37, 49, 50, 207, 218  $\overline{Y}_{\text{treatment group}} - \overline{Y}_{\text{control group}}$ , 37, 49, 138, 207, 218 Ŷ<sub>i</sub>, 102, 126, 144, 160 Ŷ, 100, 105, 112, 126, 127, 209 Z (a confounding variable), 130, 159 Z (the standard normal random variable), 173, 176, 193, 201, 202  $Z_i^X$ , 84, 94

z<sup>obs</sup>, 213, 228

*z*, 174

# Index of R and RStudio

", 11, 26 -, 9, 26 ~, 110, 128 (), 12, 26 \*, 9, 26 +, 9, 26 "13 /, 9, 26 <-, 10, 26 ; 76, 96 ==, 39, 40, 50 [], 42, 50, 61, 187, 208, 229 #, 14, 26 ##, 9 \$, 26, 40, 50, 57, 110, 128, 223 abline(), 80, 97, 111, 128 abs(), 219, 229 arguments, 12, 26 assignment operator, 10, 26 c(), 167, 195 cor(), 86, 97 CSV file, 15 data, 110, 128 data.frame(), 210, 229 dim(), 20, 26 environment, 8

Error in plot.new(): figure margins too large, 66, 109 Error: object not found, 11 exclude, 59, 96 FALSE, 11, 26, 39, 70, 96 for(i in 1:n){}, 187, 190, 194 freq, 70, 96 Google, 11 h, 80, 97 head(), 17, 26 help, 8, 60 hist(), 66, 70, 96 ifelse(), 40, 50, 56 interval, 210, 229 level, 210 lm(), 110, 115, 128, 140, 147, 161 log(), 113, 128 lty, 80, 97 margin, 65, 96 mean, 171, 175, 195 mean(), 22, 26, 96 median(), 72, 97 n, 17

238 INDEX OF R AND RSTUDIO

NA, 11, 26, 59, 96 na.omit(), 60, 96 na.rm, 60, 96 newdata, 210, 229 nrow(), 205, 229 NULL, 11, 26, 59, 96 plot(), 79, 97, 109 plots, 8, 66 pnorm(), 195 predict(), 210, 229 print(), 190, 195 prob, 181, 195 prop.table(), 58, 64, 96, 168 R console, 8 R errors, 8, 11, 13, 14, 16, 17, 66, 81, 109, 111, 187 R functions, 12, 26 R objects, 10, 26 R script, 8 read.csv(), 16, 26 replace, 181, 195 rnorm(), 171, 195 run icon, 9

sample(), 168, 195 sd, 171, 175, 195 sd(), 75, 97 setwd(), 15, 26 size, 181, 195 sqrt(), 13, 26 Stack Overflow, 11 summary()\$coef, 223, 229 t-value, 223 table(), 57, 63, 96, 168 TRUE, 11, 26, 39, 60, 96, 181, 195 v, 80, 97 var(), 76, 97 View(), 17, 26 working directory, 15, 26 workspace, 10, 14 x, 79, 97 y, 79, 97