

# Statistical Principles: PSY 200 UNM

A derivative of [openintro.org](https://openintro.org)

David M Diez

*Quantitative Analyst*

[david@openintro.org](mailto:david@openintro.org)

Christopher D Barr

*Graduate Student*

*Yale School of Management*

[chris@openintro.org](mailto:chris@openintro.org)

Mine Çetinkaya-Rundel

*Assistant Professor of the Practice*

*Department of Statistics*

*Duke University*

[mine@openintro.org](mailto:mine@openintro.org)

Instructor and adapted by:

Eric S Kruger

*Graduate Student*

*The University of New Mexico*

[esk@unm.edu](mailto:esk@unm.edu)

Copyright © 2015. Third Edition.  
Updated: August 19th, 2017.

This textbook is available under a Creative Commons license. Visit [openintro.org](http://openintro.org) for a free PDF, to download the textbook's source files, or for more information about the license.

# Contents

<b>1</b>	<b>Introduction to data</b>	<b>7</b>
1.1	Case study: using stents to prevent strokes . . . . .	7
1.2	Data basics . . . . .	9
1.3	Overview of data collection principles . . . . .	15
1.4	Observational studies and sampling strategies . . . . .	19
1.5	Experiments . . . . .	24
1.6	Examining numerical data . . . . .	26
1.7	Considering categorical data . . . . .	43
1.8	Case study: gender discrimination (special topic) . . . . .	50
1.9	Exercises . . . . .	55
<b>3</b>	<b>Distributions of random variables</b>	<b>76</b>
3.1	Normal distribution . . . . .	76
3.2	Evaluating the normal approximation . . . . .	86
3.3	Exercises . . . . .	91
<b>4</b>	<b>Foundations for inference</b>	<b>95</b>
4.1	Variability in estimates . . . . .	96
4.2	Confidence intervals . . . . .	101
4.3	Hypothesis testing . . . . .	107
4.4	Examining the Central Limit Theorem . . . . .	121
4.5	Exercises . . . . .	125
<b>5</b>	<b>Inference for numerical data</b>	<b>140</b>
5.1	One-sample means with the $t$ -distribution . . . . .	140
5.2	Paired data . . . . .	149
5.3	Difference of two means . . . . .	151
5.4	Comparing many means with ANOVA (special topic) . . . . .	160
5.5	Exercises . . . . .	172
<b>7</b>	<b>Introduction to linear regression</b>	<b>189</b>
7.1	Line fitting, residuals, and correlation . . . . .	191
7.2	Fitting a line by least squares regression . . . . .	198
7.3	Types of outliers in linear regression . . . . .	207
7.4	Inference for linear regression . . . . .	209
7.5	Exercises . . . . .	214
<b>A</b>	<b>End of chapter exercise solutions</b>	<b>230</b>

<b>B Distribution tables</b>	<b>252</b>
B.1 Normal Probability Table . . . . .	252
B.2 t-Probability Table . . . . .	255
B.3 Chi-Square Probability Table . . . . .	257

# Preface

We hope readers will take away three ideas from this book in addition to forming a foundation of statistical thinking and methods.

- (1) Statistics is an applied field with a wide range of practical applications.
- (2) You don't have to be a math guru to learn from real, interesting data.
- (3) Data are messy, and statistical tools are imperfect. But, when you understand the strengths and weaknesses of these tools, you can use them to learn about the real world.

## Textbook Modifications

This book has been modified so that it is tailored to PSY 200 Introduction to Statistics taught at The University of New Mexico Summer 2018. Please see the course syllabus if you would prefer the complete text book (it's free so why not). Specifically, chapters (2) Probability, (6) Inference for categorical data and (8) Multiple and logistic regression have been removed. Other sub-sections have been removed from each chapter. A full version of this book may be downloaded as a free PDF at [openintro.org](http://openintro.org).

## Textbook overview

The chapters of this book are as follows:

1. **Introduction to data.** Data structures, variables, summaries, graphics, and basic data collection techniques.
3. **Distributions of random variables.** Introduction to the normal model and other key distributions.
4. **Foundations for inference.** General ideas for statistical inference in the context of estimating the population mean.
5. **Inference for numerical data.** Inference for one or two sample means using the  $t$ -distribution, and also comparisons of many means using ANOVA.
7. **Introduction to linear regression.** An introduction to regression with two variables. Most of this chapter could be covered after Chapter 1.

*OpenIntro Statistics* was written to allow flexibility in choosing and ordering course topics. The material is divided into two pieces: main text and special topics. The main text has been structured to bring statistical inference and modeling closer to the front of a course. Special topics, labeled in the table of contents and in section titles, may be added to a course as they arise naturally in the curriculum.

## Videos for sections and calculators

The  icon indicates that a section or topic has a video overview readily available. The icons are hyperlinked in the textbook PDF, and the videos may also be found at

[www.openintro.org/stat/videos.php](http://www.openintro.org/stat/videos.php)

There are also slides that provide further explanation of via case studies. To access in the book click on the  to open the slides in Google Slides.

## Examples, exercises, and appendices

Examples and Guided Practice throughout the textbook may be identified by their distinctive bullets:

- **Example 0.1** Large filled bullets signal the start of an example.  
Full solutions to examples are provided and may include an accompanying table or figure.
- **Guided Practice 0.2** Large empty bullets signal to readers that an exercise has been inserted into the text for additional practice and guidance. Students may find it useful to fill in the bullet after understanding or successfully completing the exercise. Solutions are provided for all Guided Practice in footnotes.<sup>1</sup>

There are exercises at the end of each chapter for practice or homework assignments. Odd-numbered exercise solutions are in Appendix A. Probability tables for the normal,  $t$ , and chi-square distributions are in Appendix B.

## OpenIntro, online resources, and getting involved

OpenIntro is an organization focused on developing free and affordable education materials. *OpenIntro Statistics* is intended for introductory statistics courses at the college level. We offer another title, *Advanced High School Statistics*, for high school courses.

We encourage anyone learning or teaching statistics to visit [openintro.org](http://openintro.org) and get involved. We also provide many free online resources, including free course software. Data sets for this textbook are available on the website and through a companion R package.<sup>2</sup> All of these resources are free and may be used with or without this textbook as a companion.

We value your feedback. If there is a particular component of the project you especially like or think needs improvement, we want to hear from you. You may find our contact information on the title page of this book or on the *About* section of [openintro.org](http://openintro.org).

## Acknowledgements

This project would not be possible without the passion and dedication of all those involved. The authors would like to thank the OpenIntro Staff for their involvement and ongoing contributions. We are also very grateful to the hundreds of students and instructors who have provided us with valuable feedback over the last several years.

---

<sup>1</sup>Full solutions are located down here in the footnote!

<sup>2</sup>Diez DM, Barr CD, Çetinkaya-Rundel M. 2015. `openintro`: OpenIntro data sets and supplement functions. [github.com/OpenIntroOrg/openintro-r-package](https://github.com/OpenIntroOrg/openintro-r-package).

# Chapter 1

## Introduction to data

Scientists seek to answer questions using rigorous methods and careful observations. These observations – collected from the likes of field notes, surveys, and experiments – form the backbone of a statistical investigation and are called **data**. Statistics is the study of how best to collect, analyze, and draw conclusions from data. It is helpful to put statistics in the context of a general process of investigation:

1. Identify a question or problem.
2. Collect relevant data on the topic.
3. Analyze the data.
4. Form a conclusion.

Statistics as a subject focuses on making stages 2-4 objective, rigorous, and efficient. That is, statistics has three primary components: How best can we collect data? How should it be analyzed? And what can we infer from the analysis?

The topics scientists investigate are as diverse as the questions they ask. However, many of these investigations can be addressed with a small number of data collection techniques, analytic tools, and fundamental concepts in statistical inference. This chapter provides a glimpse into these and other themes we will encounter throughout the rest of the book. We introduce the basic principles of each branch and learn some tools along the way. We will encounter applications from other fields, some of which are not typically associated with science but nonetheless can benefit from statistical study.

### 1.1 Case study: using stents to prevent strokes

Section 1.1 introduces a classic challenge in statistics: evaluating the efficacy of a medical treatment. Terms in this section, and indeed much of this chapter, will all be revisited later in the text. The plan for now is simply to get a sense of the role statistics can play in practice.

In this section we will consider an experiment that studies effectiveness of stents in treating patients at risk of stroke.<sup>1</sup> Stents are devices put inside blood vessels that assist

---

<sup>1</sup>Chimowitz MI, Lynn MJ, Derdeyn CP, et al. 2011. Stenting versus Aggressive Medical Therapy for Intracranial Arterial Stenosis. New England Journal of Medicine 365:993-1003. [www.nejm.org/doi/full/10.1056/NEJMoa1105335](http://www.nejm.org/doi/full/10.1056/NEJMoa1105335). NY Times article reporting on the study: [www.nytimes.com/2011/09/08/health/research/08stent.html](http://www.nytimes.com/2011/09/08/health/research/08stent.html).

in patient recovery after cardiac events and reduce the risk of an additional heart attack or death. Many doctors have hoped that there would be similar benefits for patients at risk of stroke. We start by writing the principal question the researchers hope to answer:

Does the use of stents reduce the risk of stroke?

The researchers who asked this question collected data on 451 at-risk patients. Each volunteer patient was randomly assigned to one of two groups:

**Treatment group.** Patients in the treatment group received a stent and medical management. The medical management included medications, management of risk factors, and help in lifestyle modification.

**Control group.** Patients in the control group received the same medical management as the treatment group, but they did not receive stents.

Researchers randomly assigned 224 patients to the treatment group and 227 to the control group. In this study, the control group provides a reference point against which we can measure the medical impact of stents in the treatment group.

Researchers studied the effect of stents at two time points: 30 days after enrollment and 365 days after enrollment. The results of 5 patients are summarized in Table 1.1. Patient outcomes are recorded as “stroke” or “no event”, representing whether or not the patient had a stroke at the end of a time period.

Patient	group	0-30 days	0-365 days
1	treatment	no event	no event
2	treatment	stroke	stroke
3	treatment	no event	no event
:	:	:	
450	control	no event	no event
451	control	no event	no event

Table 1.1: Results for five patients from the stent study.

Considering data from each patient individually would be a long, cumbersome path towards answering the original research question. Instead, performing a statistical data analysis allows us to consider all of the data at once. Table 1.2 summarizes the raw data in a more helpful way. In this table, we can quickly see what happened over the entire study. For instance, to identify the number of patients in the treatment group who had a stroke within 30 days, we look on the left-side of the table at the intersection of the treatment and stroke: 33.

	0-30 days		0-365 days	
	stroke	no event	stroke	no event
treatment	33	191	45	179
control	13	214	28	199
Total	46	405	73	378

Table 1.2: Descriptive statistics for the stent study.

Ⓐ **Guided Practice 1.1** Of the 224 patients in the treatment group, 45 had a stroke by the end of the first year. Using these two numbers, compute the proportion of patients in the treatment group who had a stroke by the end of their first year. (Please note: answers to all Guided Practice exercises are provided using footnotes.)<sup>2</sup>

We can compute summary statistics from the table. A **summary statistic** is a single number summarizing a large amount of data.<sup>3</sup> For instance, the primary results of the study after 1 year could be described by two summary statistics: the proportion of people who had a stroke in the treatment and control groups.

Proportion who had a stroke in the treatment (stent) group:  $45/224 = 0.20 = 20\%$ .

Proportion who had a stroke in the control group:  $28/227 = 0.12 = 12\%$ .

These two summary statistics are useful in looking for differences in the groups, and we are in for a surprise: an additional 8% of patients in the treatment group had a stroke! This is important for two reasons. First, it is contrary to what doctors expected, which was that stents would *reduce* the rate of strokes. Second, it leads to a statistical question: do the data show a “real” difference between the groups?

This second question is subtle. Suppose you flip a coin 100 times. While the chance a coin lands heads in any given coin flip is 50%, we probably won’t observe exactly 50 heads. This type of fluctuation is part of almost any type of data generating process. It is possible that the 8% difference in the stent study is due to this natural variation. However, the larger the difference we observe (for a particular sample size), the less believable it is that the difference is due to chance. So what we are really asking is the following: is the difference so large that we should reject the notion that it was due to chance?

While we don’t yet have our statistical tools to fully address this question on our own, we can comprehend the conclusions of the published analysis: there was compelling evidence of harm by stents in this study of stroke patients.

**Be careful:** do not generalize the results of this study to all patients and all stents. This study looked at patients with very specific characteristics who volunteered to be a part of this study and who may not be representative of all stroke patients. In addition, there are many types of stents and this study only considered the self-expanding Wingspan stent (Boston Scientific). However, this study does leave us with an important lesson: we should keep our eyes open for surprises.

## 1.2 Data basics

Effective presentation and description of data is a first step in most analyses. This section introduces one structure for organizing data as well as some terminology that will be used throughout this book.

### 1.2.1 Observations, variables, and data matrices

Table 1.3 displays rows 1, 2, 3, and 50 of a data set concerning 50 emails received during early 2012. These observations will be referred to as the `email50` data set, and they are a random sample from a larger data set that we will see in Section 1.7.

---

<sup>2</sup>The proportion of the 224 patients who had a stroke within 365 days:  $45/224 = 0.20$ .

<sup>3</sup>Formally, a summary statistic is a value computed from the data. Some summary statistics are more useful than others.

	spam	num_char	line_breaks	format	number
1	no	21,705	551	html	small
2	no	7,011	183	html	big
3	yes	631	28	text	none
:	:	:	:	:	:
50	no	15,829	242	html	small

Table 1.3: Four rows from the `email150` data matrix.

variable	description
<code>spam</code>	Specifies whether the message was spam
<code>num_char</code>	The number of characters in the email
<code>line_breaks</code>	The number of line breaks in the email (not including text wrapping)
<code>format</code>	Indicates if the email contained special formatting, such as bolding, tables, or links, which would indicate the message is in HTML format
<code>number</code>	Indicates whether the email contained no number, a small number (under 1 million), or a large number

Table 1.4: Variables and their descriptions for the `email150` data set.

Each row in the table represents a single email or **case**.<sup>4</sup> The columns represent characteristics, called **variables**, for each of the emails. For example, the first row represents email 1, which is a not spam, contains 21,705 characters, 551 line breaks, is written in HTML format, and contains only small numbers.

In practice, it is especially important to ask clarifying questions to ensure important aspects of the data are understood. For instance, it is always important to be sure we know what each variable means and the units of measurement. Descriptions of all five email variables are given in Table 1.4.

The data in Table 1.3 represent a **data matrix**, which is a common way to organize data. Each row of a data matrix corresponds to a unique case, and each column corresponds to a variable. A data matrix for the stroke study introduced in Section 1.1 is shown in Table 1.1 on page 8, where the cases were patients and there were three variables recorded for each patient.

Data matrices are a convenient way to record and store data. If another individual or case is added to the data set, an additional row can be easily added. Similarly, another column can be added for a new variable.

- **Guided Practice 1.2** We consider a publicly available data set that summarizes information about the 3,143 counties in the United States, and we call this the `county` data set. This data set includes information about each county: its name, the state where it resides, its population in 2000 and 2010, per capita federal spending, poverty rate, and five additional characteristics. How might these data be organized in a data matrix? Reminder: look in the footnotes for answers to in-text exercises.<sup>5</sup>

Seven rows of the `county` data set are shown in Table 1.5, and the variables are summarized in Table 1.6. These data were collected from the US Census website.<sup>6</sup>

<sup>4</sup>A case is also sometimes called a **unit of observation** or an **observational unit**.

<sup>5</sup>Each county may be viewed as a case, and there are eleven pieces of information recorded for each case. A table with 3,143 rows and 11 columns could hold these data, where each row represents a county and each column represents a particular piece of information.

<sup>6</sup>[quickfacts.census.gov/qfd/index.html](http://quickfacts.census.gov/qfd/index.html)

	<b>name</b>	<b>state</b>	<b>pop2000</b>	<b>pop2010</b>	<b>fed_spend</b>	<b>poverty</b>	<b>homeownership</b>	<b>multiunit</b>	<b>income</b>	<b>med_income</b>	<b>smoking_ban</b>
1	Autauga	AL	43671	54571	6.068	10.6	77.5	7.2	24568	53255	none
2	Baldwin	AL	140415	182265	6.140	12.2	76.7	22.6	26469	50147	none
3	Barbour	AL	29038	27457	8.752	25.0	68.0	11.1	15875	33219	none
4	Bibb	AL	20826	22915	7.122	12.6	82.9	6.6	19918	41770	none
5	Blount	AL	51024	57322	5.131	13.4	82.0	3.7	21070	45549	none
:	:	:	:	:	:	:	:	:	:	:	:
3142	Washakie	WY	8289	8533	8.714	5.6	70.9	10.0	28557	48379	none
3143	Weston	WY	6644	7208	6.695	7.9	77.9	6.5	28463	53853	none

Table 1.5: Seven rows from the county data set.

<b>variable</b>	<b>description</b>
<b>name</b>	County name
<b>state</b>	State where the county resides (also including the District of Columbia)
<b>pop2000</b>	Population in 2000
<b>pop2010</b>	Population in 2010
<b>fed_spend</b>	Federal spending per capita
<b>poverty</b>	Percent of the population in poverty
<b>homeownership</b>	Percent of the population that lives in their own home or lives with the owner (e.g. children living with parents who own the home)
<b>multiunit</b>	Percent of living units that are in multi-unit structures (e.g. apartments)
<b>income</b>	Income per capita
<b>med_income</b>	Median household income for the county, where a household's income equals the total income of its occupants who are 15 years or older
<b>smoking_ban</b>	Type of county-wide smoking ban in place at the end of 2011, which takes one of three values: <b>none</b> , <b>partial</b> , or <b>comprehensive</b> , where a <b>comprehensive</b> ban means smoking was not permitted in restaurants, bars, or workplaces, and <b>partial</b> means smoking was banned in at least one of those three locations

Table 1.6: Variables and their descriptions for the county data set.



Figure 1.7: Breakdown of variables into their respective types.

### 1.2.2 Types of variables

Examine the `fed_spend`, `pop2010`, `state`, and `smoking_ban` variables in the `county` data set. Each of these variables is inherently different from the other three yet many of them share certain characteristics.

First consider `fed_spend`, which is said to be a **numerical** variable since it can take a wide range of numerical values, and it is sensible to add, subtract, or take averages with those values. On the other hand, we would not classify a variable reporting telephone area codes as numerical since their average, sum, and difference have no clear meaning.

The `pop2010` variable is also numerical, although it seems to be a little different than `fed_spend`. This variable of the population count can only take whole non-negative numbers (0, 1, 2, ...). For this reason, the population variable is said to be **discrete** since it can only take numerical values with jumps. On the other hand, the federal spending variable is said to be **continuous**.

The variable `state` can take up to 51 values after accounting for Washington, DC: `AL`, ..., and `WY`. Because the responses themselves are categories, `state` is called a **categorical** variable, and the possible values are called the variable's **levels**.

Finally, consider the `smoking_ban` variable, which describes the type of county-wide smoking ban and takes values `none`, `partial`, or `comprehensive` in each county. This variable seems to be a hybrid: it is a categorical variable but the levels have a natural ordering. A variable with these properties is called an **ordinal** variable, while a regular categorical variable without this type of special ordering is called a **nominal** variable. To simplify analyses, any ordinal variables in this book will be treated as categorical variables.

- **Example 1.3** Data were collected about students in a statistics course. Three variables were recorded for each student: number of siblings, student height, and whether the student had previously taken a statistics course. Classify each of the variables as continuous numerical, discrete numerical, or categorical.

---

The number of siblings and student height represent numerical variables. Because the number of siblings is a count, it is discrete. Height varies continuously, so it is a continuous numerical variable. The last variable classifies students into two categories – those who have and those who have not taken a statistics course – which makes this variable categorical.

- **Guided Practice 1.4** Consider the variables `group` and `outcome` (at 30 days) from the stent study in Section 1.1. Are these numerical or categorical variables?<sup>7</sup>

---

<sup>7</sup>There are only two possible values for each variable, and in both cases they describe categories. Thus, each is a categorical variable.

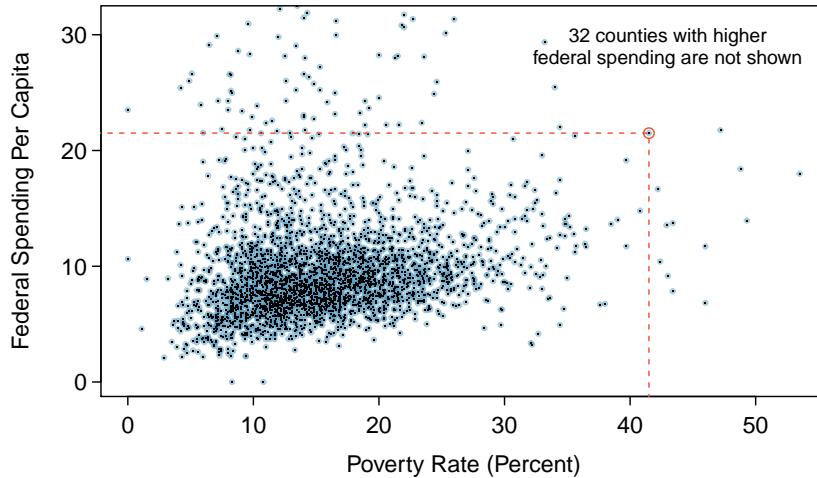


Figure 1.8: A scatterplot showing `fed_spend` against `poverty`. Owsley County of Kentucky, with a poverty rate of 41.5% and federal spending of \$21.50 per capita, is highlighted.

### 1.2.3 Relationships between variables

Many analyses are motivated by a researcher looking for a relationship between two or more variables. A social scientist may like to answer some of the following questions:

- (1) Is federal spending, on average, higher or lower in counties with high rates of poverty?
- (2) If homeownership is lower than the national average in one county, will the percent of multi-unit structures in that county likely be above or below the national average?
- (3) Which counties have a higher average income: those that enact one or more smoking bans or those that do not?

To answer these questions, data must be collected, such as the `county` data set shown in Table 1.5. Examining summary statistics could provide insights for each of the three questions about counties. Additionally, graphs can be used to visually summarize data and are useful for answering such questions as well.

Scatterplots are one type of graph used to study the relationship between two numerical variables. Figure 1.8 compares the variables `fed_spend` and `poverty`. Each point on the plot represents a single county. For instance, the highlighted dot corresponds to County 1088 in the `county` data set: Owsley County, Kentucky, which had a poverty rate of 41.5% and federal spending of \$21.50 per capita. The scatterplot suggests a relationship between the two variables: counties with a high poverty rate also tend to have slightly more federal spending. We might brainstorm as to why this relationship exists and investigate each idea to determine which is the most reasonable explanation.

-  **Guided Practice 1.5** Examine the variables in the `email150` data set, which are described in Table 1.4 on page 10. Create two questions about the relationships between these variables that are of interest to you.<sup>8</sup>

---

<sup>8</sup>Two sample questions: (1) Intuition suggests that if there are many line breaks in an email then there also would tend to be many characters: does this hold true? (2) Is there a connection between whether an email format is plain text (versus HTML) and whether it is a spam message?

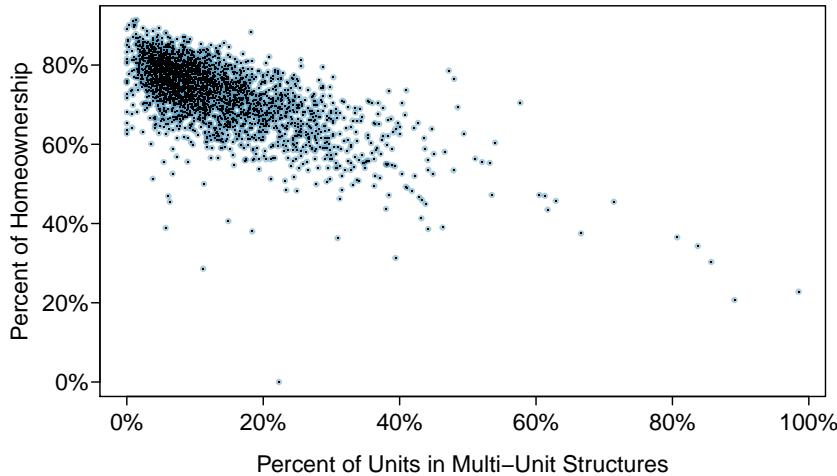


Figure 1.9: A scatterplot of homeownership versus the percent of units that are in multi-unit structures for all 3,143 counties. Interested readers may find an image of this plot with an additional third variable, county population, presented at [www.openintro.org/stat/down/MHP.png](http://www.openintro.org/stat/down/MHP.png).

The `fed_spend` and `poverty` variables are said to be associated because the plot shows a discernible pattern. When two variables show some connection with one another, they are called **associated** variables. Associated variables can also be called **dependent** variables and vice-versa.

- **Example 1.6** This example examines the relationship between homeownership and the percent of units in multi-unit structures (e.g. apartments, condos), which is visualized using a scatterplot in Figure 1.9. Are these variables associated?

It appears that the larger the fraction of units in multi-unit structures, the lower the homeownership rate. Since there is some relationship between the variables, they are associated.

Because there is a downward trend in Figure 1.9 – counties with more units in multi-unit structures are associated with lower homeownership – these variables are said to be **negatively associated**. A **positive association** is shown in the relationship between the `poverty` and `fed_spend` variables represented in Figure 1.8, where counties with higher poverty rates tend to receive more federal spending per capita.

If two variables are not associated, then they are said to be **independent**. That is, two variables are independent if there is no evident relationship between the two.

#### Associated or independent, not both

A pair of variables are either related in some way (associated) or not (independent). No pair of variables is both associated and independent.

## 1.3 Overview of data collection principles

The first step in conducting research is to identify topics or questions that are to be investigated. A clearly laid out research question is helpful in identifying what subjects or cases should be studied and what variables are important. It is also important to consider *how* data are collected so that they are reliable and help achieve the research goals.

### 1.3.1 Populations and samples

Consider the following three research questions:

1. What is the average mercury content in swordfish in the Atlantic Ocean?
2. Over the last 5 years, what is the average time to complete a degree for Duke undergraduate students?
3. Does a new drug reduce the number of deaths in patients with severe heart disease?

Each research question refers to a target **population**. In the first question, the target population is all swordfish in the Atlantic ocean, and each fish represents a case. Often times, it is too expensive to collect data for every case in a population. Instead, a sample is taken. A **sample** represents a subset of the cases and is often a small fraction of the population. For instance, 60 swordfish (or some other number) in the population might be selected, and this sample data may be used to provide an estimate of the population average and answer the research question.

 **Guided Practice 1.7** For the second and third questions above, identify the target population and what represents an individual case.<sup>9</sup>

### 1.3.2 Anecdotal evidence

Consider the following possible responses to the three research questions:

1. A man on the news got mercury poisoning from eating swordfish, so the average mercury concentration in swordfish must be dangerously high.
2. I met two students who took more than 7 years to graduate from Duke, so it must take longer to graduate at Duke than at many other colleges.
3. My friend's dad had a heart attack and died after they gave him a new heart disease drug, so the drug must not work.

Each conclusion is based on data. However, there are two problems. First, the data only represent one or two cases. Second, and more importantly, it is unclear whether these cases are actually representative of the population. Data collected in this haphazard fashion are called **anecdotal evidence**.

---

<sup>9</sup>(2) Notice that the first question is only relevant to students who complete their degree; the average cannot be computed using a student who never finished her degree. Thus, only Duke undergraduate students who have graduated in the last five years represent cases in the population under consideration. Each such student would represent an individual case. (3) A person with severe heart disease represents a case. The population includes all people with severe heart disease.



Figure 1.10: In February 2010, some media pundits cited one large snow storm as valid evidence against global warming. As comedian Jon Stewart pointed out, “It’s one storm, in one region, of one country.”

### Anecdotal evidence

Be careful of data collected in a haphazard fashion. Such evidence may be true and verifiable, but it may only represent extraordinary cases.

Anecdotal evidence typically is composed of unusual cases that we recall based on their striking characteristics. For instance, we are more likely to remember the two people we met who took 7 years to graduate than the six others who graduated in four years. Instead of looking at the most unusual cases, we should examine a sample of many cases that represent the population.

### 1.3.3 Sampling from a population

We might try to estimate the time to graduation for Duke undergraduates in the last 5 years by collecting a sample of students. All graduates in the last 5 years represent the *population*, and graduates who are selected for review are collectively called the *sample*. In general, we always seek to *randomly* select a sample from a population. The most basic type of random selection is equivalent to how raffles are conducted. For example, in selecting graduates, we could write each graduate’s name on a raffle ticket and draw 100 tickets. The selected names would represent a random sample of 100 graduates.

Why pick a sample randomly? Why not just pick a sample by hand? Consider the following scenario.

- **Example 1.8** Suppose we ask a student who happens to be majoring in nutrition to select several graduates for the study. What kind of students do you think she might collect? Do you think her sample would be representative of all graduates?

Perhaps she would pick a disproportionate number of graduates from health-related fields. Or perhaps her selection would be well-representative of the population. When selecting samples by hand, we run the risk of picking a *biased* sample, even if that bias is unintentional or difficult to discern.

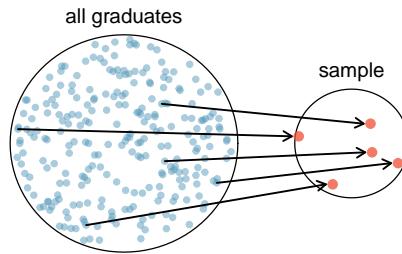


Figure 1.11: In this graphic, five graduates are randomly selected from the population to be included in the sample.

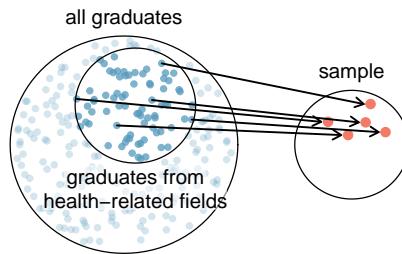


Figure 1.12: Instead of sampling from all graduates equally, a nutrition major might inadvertently pick graduates with health-related majors disproportionately often.

If someone was permitted to pick and choose exactly which graduates were included in the sample, it is entirely possible that the sample could be skewed to that person's interests, which may be entirely unintentional. This introduces **bias** into a sample. Sampling randomly helps resolve this problem. The most basic random sample is called a **simple random sample**, and which is equivalent to using a raffle to select cases. This means that each case in the population has an equal chance of being included and there is no implied connection between the cases in the sample.

The act of taking a simple random sample helps minimize bias. However, bias can crop up in other ways. Even when people are picked at random, e.g. for surveys, caution must be exercised if the **non-response** is high. For instance, if only 30% of the people randomly sampled for a survey actually respond, then it is unclear whether the results are **representative** of the entire population. This **non-response bias** can skew results.

Another common downfall is a **convenience sample**, where individuals who are easily accessible are more likely to be included in the sample. For instance, if a political survey is done by stopping people walking in the Bronx, this will not represent all of New York City. It is often difficult to discern what sub-population a convenience sample represents.

 **Guided Practice 1.9** We can easily access ratings for products, sellers, and companies through websites. These ratings are based only on those people who go out of their way to provide a rating. If 50% of online reviews for a product are negative, do you think this means that 50% of buyers are dissatisfied with the product?<sup>10</sup>

<sup>10</sup>Answers will vary. From our own anecdotal experiences, we believe people tend to rant more about products that fell below expectations than rave about those that perform as expected. For this reason, we suspect there is a negative bias in product ratings on sites like Amazon. However, since our experiences may not be representative, we also keep an open mind.

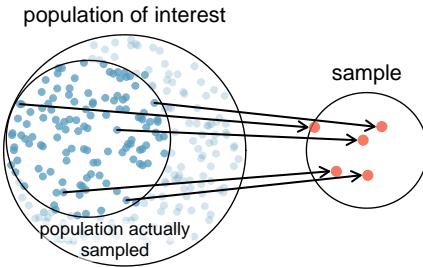


Figure 1.13: Due to the possibility of non-response, surveys studies may only reach a certain group within the population. It is difficult, and often times impossible, to completely fix this problem.

### 1.3.4 Explanatory and response variables

Consider the following question from page 13 for the `county` data set:

- (1) Is federal spending, on average, higher or lower in counties with high rates of poverty?

If we suspect poverty might affect spending in a county, then poverty is the **explanatory** variable and federal spending is the **response** variable in the relationship.<sup>11</sup> If there are many variables, it may be possible to consider a number of them as explanatory variables.

**TIP: Explanatory and response variables**

To identify the explanatory variable in a pair of variables, identify which of the two is suspected of affecting the other and plan an appropriate analysis.

explanatory variable	$\xrightarrow{\text{might affect}}$	response variable
-------------------------	-------------------------------------	----------------------

**Caution: association does not imply causation**

Labeling variables as *explanatory* and *response* does not guarantee the relationship between the two is actually causal, even if there is an association identified between the two variables. We use these labels only to keep track of which variable we suspect affects the other.

In some cases, there is no explanatory or response variable. Consider the following question from page 13:

- (2) If homeownership is lower than the national average in one county, will the percent of multi-unit structures in that county likely be above or below the national average?

It is difficult to decide which of these variables should be considered the explanatory and response variable, i.e. the direction is ambiguous, so no explanatory or response labels are suggested here.

<sup>11</sup>Sometimes the explanatory variable is called the **independent** variable and the response variable is called the **dependent** variable. However, this becomes confusing since a *pair* of variables might be independent or dependent, so we avoid this language.

### 1.3.5 Introducing observational studies and experiments

There are two primary types of data collection: observational studies and experiments.

Researchers perform an **observational study** when they collect data in a way that does not directly interfere with how the data arise. For instance, researchers may collect information via surveys, review medical or company records, or follow a **cohort** of many similar individuals to study why certain diseases might develop. In each of these situations, researchers merely observe the data that arise. In general, observational studies can provide evidence of a naturally occurring association between variables, but they cannot by themselves show a causal connection.

When researchers want to investigate the possibility of a causal connection, they conduct an **experiment**. Usually there will be both an explanatory and a response variable. For instance, we may suspect administering a drug will reduce mortality in heart attack patients over the following year. To check if there really is a causal connection between the explanatory variable and the response, researchers will collect a sample of individuals and split them into groups. The individuals in each group are *assigned* a treatment. When individuals are randomly assigned to a group, the experiment is called a **randomized experiment**. For example, each heart attack patient in the drug trial could be randomly assigned, perhaps by flipping a coin, into one of two groups: the first group receives a **placebo** (fake treatment) and the second group receives the drug. See the case study in Section 1.1 for another example of an experiment, though that study did not employ a placebo.

#### TIP: association $\neq$ causation

In general, association does not imply causation, and causation can only be inferred from a randomized experiment.

## 1.4 Observational studies and sampling strategies

### 1.4.1 Observational studies

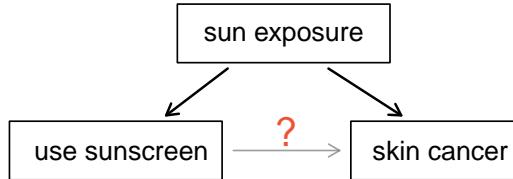
Generally, data in observational studies are collected only by monitoring what occurs, while experiments require the primary explanatory variable in a study be assigned for each subject by the researchers.

Making causal conclusions based on experiments is often reasonable. However, making the same causal conclusions based on observational data can be treacherous and is not recommended. Thus, observational studies are generally only sufficient to show associations.

- Ⓐ **Guided Practice 1.10** Suppose an observational study tracked sunscreen use and skin cancer, and it was found that the more sunscreen someone used, the more likely the person was to have skin cancer. Does this mean sunscreen *causes* skin cancer?<sup>12</sup>

Some previous research tells us that using sunscreen actually reduces skin cancer risk, so maybe there is another variable that can explain this hypothetical association between sunscreen usage and skin cancer. One important piece of information that is absent is sun exposure. If someone is out in the sun all day, she is more likely to use sunscreen *and* more likely to get skin cancer. Exposure to the sun is unaccounted for in the simple investigation.

<sup>12</sup>No. See the paragraph following the exercise for an explanation.



Sun exposure is what is called a **confounding variable**,<sup>13</sup> which is a variable that is correlated with both the explanatory and response variables. While one method to justify making causal conclusions from observational studies is to exhaust the search for confounding variables, there is no guarantee that all confounding variables can be examined or measured.

In the same way, the `county` data set is an observational study with confounding variables, and its data cannot easily be used to make causal conclusions.

- Ⓐ **Guided Practice 1.11** Figure 1.9 shows a negative association between the homeownership rate and the percentage of multi-unit structures in a county. However, it is unreasonable to conclude that there is a causal relationship between the two variables. Suggest one or more other variables that might explain the relationship visible in Figure 1.9.<sup>14</sup>

Observational studies come in two forms: prospective and retrospective studies. A **prospective study** identifies individuals and collects information as events unfold. For instance, medical researchers may identify and follow a group of similar individuals over many years to assess the possible influences of behavior on cancer risk. One example of such a study is The Nurses' Health Study, started in 1976 and expanded in 1989.<sup>15</sup> This prospective study recruits registered nurses and then collects data from them using questionnaires. **Retrospective studies** collect data after events have taken place, e.g. researchers may review past events in medical records. Some data sets, such as `county`, may contain both prospectively- and retrospectively-collected variables. Local governments prospectively collect some variables as events unfolded (e.g. retail sales) while the federal government retrospectively collected others during the 2010 census (e.g. county population counts).

### 1.4.2 Four sampling methods (special topic)

Almost all statistical methods are based on the notion of implied randomness. If observational data are not collected in a random framework from a population, these statistical methods – the estimates and errors associated with the estimates – are not reliable. Here we consider four random sampling techniques: simple, stratified, cluster, and multistage sampling. Figures 1.14 and 1.15 provide graphical representations of these techniques.

**Simple random sampling** is probably the most intuitive form of random sampling. Consider the salaries of Major League Baseball (MLB) players, where each player is a member of one of the league's 30 teams. To take a simple random sample of 120 baseball players and their salaries from the 2010 season, we could write the names of that season's

<sup>13</sup>Also called a **lurking variable**, **confounding factor**, or a **confounder**.

<sup>14</sup>Answers will vary. Population density may be important. If a county is very dense, then this may require a larger fraction of residents to live in multi-unit structures. Additionally, the high density may contribute to increases in property value, making homeownership infeasible for many residents.

<sup>15</sup>[www.channing.harvard.edu/nhs](http://www.channing.harvard.edu/nhs)

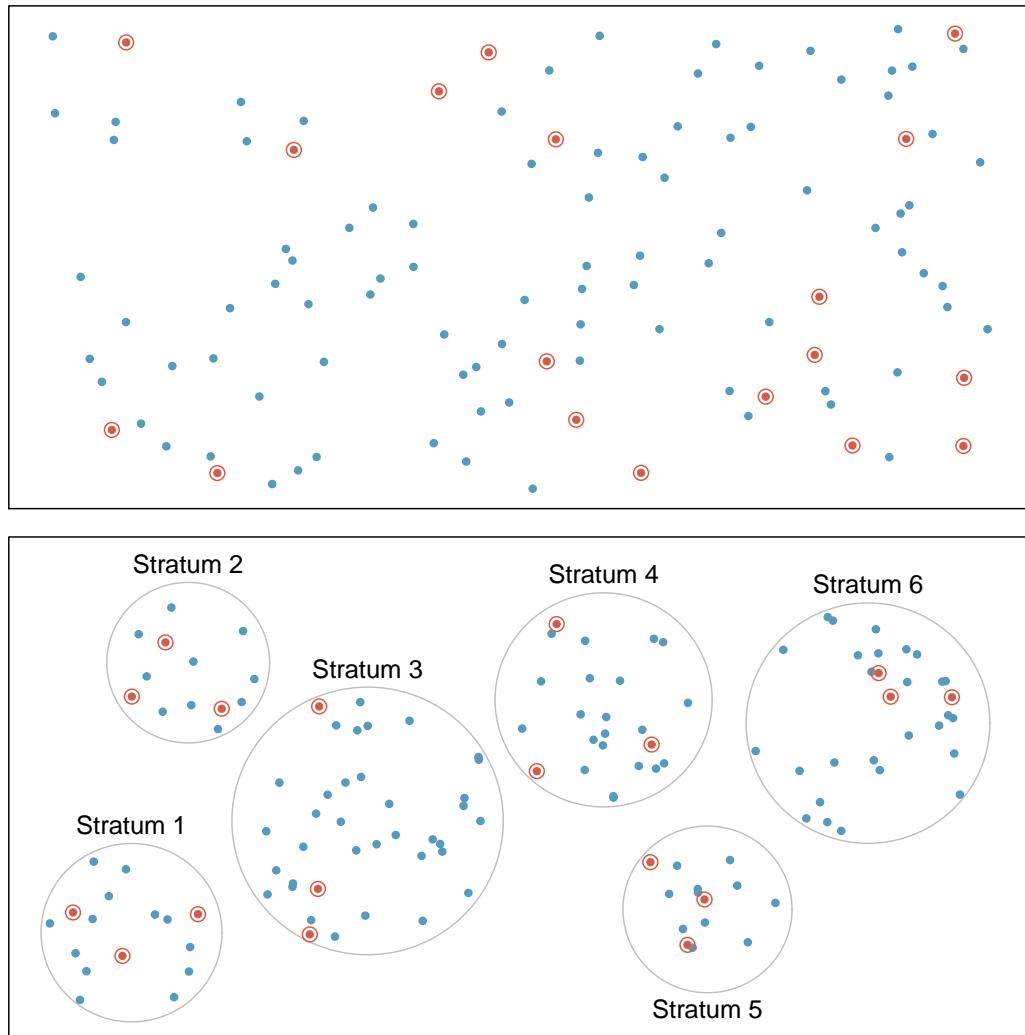


Figure 1.14: Examples of simple random and stratified sampling. In the top panel, simple random sampling was used to randomly select the 18 cases. In the bottom panel, stratified sampling was used: cases were grouped into strata, then simple random sampling was employed within each stratum.

828 players onto slips of paper, drop the slips into a bucket, shake the bucket around until we are sure the names are all mixed up, then draw out slips until we have the sample of 120 players. In general, a sample is referred to as “simple random” if each case in the population has an equal chance of being included in the final sample *and* knowing that a case is included in a sample does not provide useful information about which other cases are included.

**Stratified sampling** is a divide-and-conquer sampling strategy. The population is divided into groups called **strata**. The strata are chosen so that similar cases are grouped together, then a second sampling method, usually simple random sampling, is employed within each stratum. In the baseball salary example, the teams could represent the strata, since some teams have a lot more money (up to 4 times as much!). Then we might randomly sample 4 players from each team for a total of 120 players.

Stratified sampling is especially useful when the cases in each stratum are very similar with respect to the outcome of interest. The downside is that analyzing data from a stratified sample is a more complex task than analyzing data from a simple random sample. The analysis methods introduced in this book would need to be extended to analyze data collected using stratified sampling.

- **Example 1.12** Why would it be good for cases within each stratum to be very similar?

---

We might get a more stable estimate for the subpopulation in a stratum if the cases are very similar. These improved estimates for each subpopulation will help us build a reliable estimate for the full population.

In a **cluster sample**, we break up the population into many groups, called **clusters**. Then we sample a fixed number of clusters and include all observations from each of those clusters in the sample. A **multistage sample** is like a cluster sample, but rather than keeping all observations in each cluster, we collect a random sample within each selected cluster.

Sometimes cluster or multistage sampling can be more economical than the alternative sampling techniques. Also, unlike stratified sampling, these approaches are most helpful when there is a lot of case-to-case variability within a cluster but the clusters themselves don't look very different from one another. For example, if neighborhoods represented clusters, then cluster or multistage sampling work best when the neighborhoods are very diverse. A downside of these methods is that more advanced analysis techniques are typically required, though the methods in this book can be extended to handle such data.

- **Example 1.13** Suppose we are interested in estimating the malaria rate in a densely tropical portion of rural Indonesia. We learn that there are 30 villages in that part of the Indonesian jungle, each more or less similar to the next. Our goal is to test 150 individuals for malaria. What sampling method should be employed?

---

A simple random sample would likely draw individuals from all 30 villages, which could make data collection extremely expensive. Stratified sampling would be a challenge since it is unclear how we would build strata of similar individuals. However, cluster sampling or multistage sampling seem like very good ideas. If we decided to use multistage sampling, we might randomly select half of the villages, then randomly select 10 people from each. This would probably reduce our data collection costs substantially in comparison to a simple random sample, and this approach would still give us reliable information.

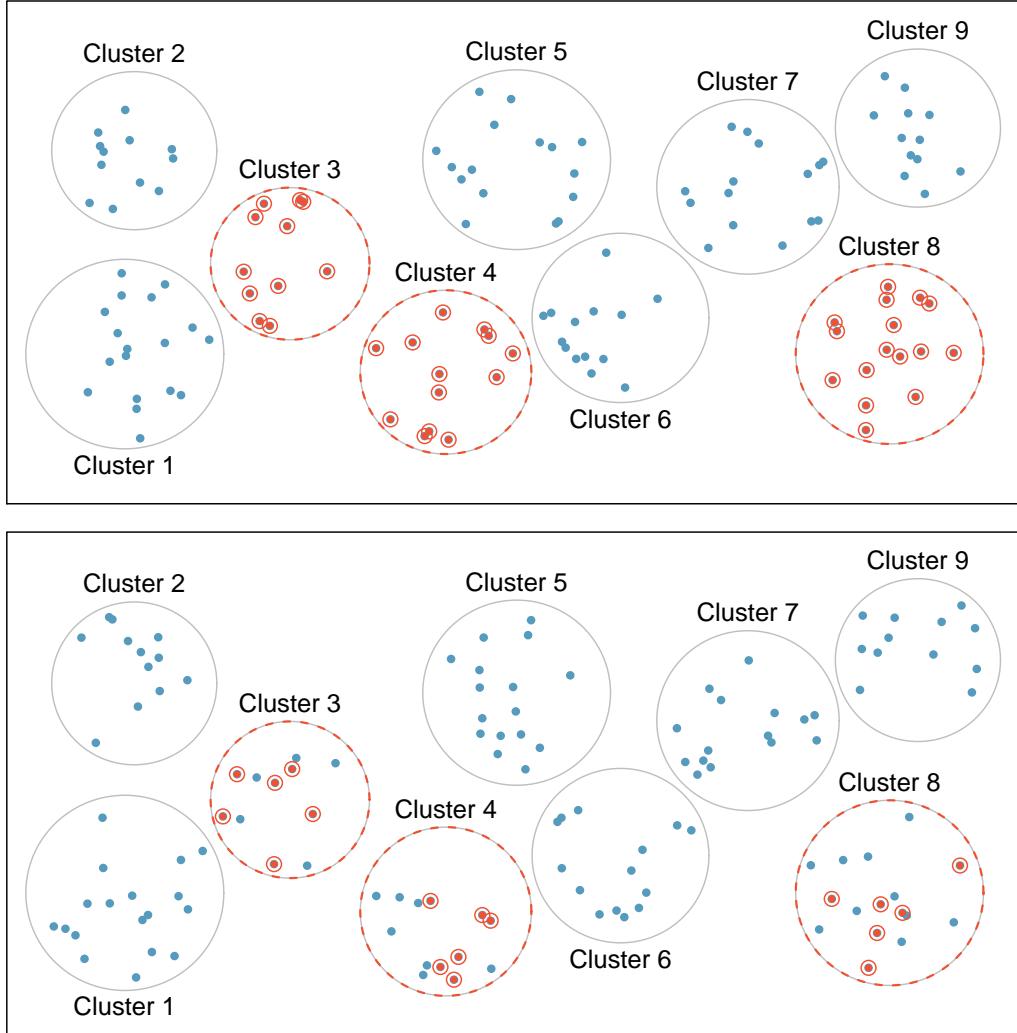


Figure 1.15: Examples of cluster and multistage sampling. In the top panel, cluster sampling was used. Here, data were binned into nine clusters, three of these clusters were sampled, and all observations within these three cluster were included in the sample. In the bottom panel, multistage sampling was used. It differs from cluster sampling in that of the clusters selected, we randomly select a subset of each cluster to be included in the sample.

## 1.5 Experiments

Studies where the researchers assign treatments to cases are called **experiments**. When this assignment includes randomization, e.g. using a coin flip to decide which treatment a patient receives, it is called a **randomized experiment**. Randomized experiments are fundamentally important when trying to show a causal connection between two variables.

### 1.5.1 Principles of experimental design

Randomized experiments are generally built on four principles.

**Controlling.** Researchers assign treatments to cases, and they do their best to **control** any other differences in the groups. For example, when patients take a drug in pill form, some patients take the pill with only a sip of water while others may have it with an entire glass of water. To control for the effect of water consumption, a doctor may ask all patients to drink a 12 ounce glass of water with the pill.

**Randomization.** Researchers randomize patients into treatment groups to account for variables that cannot be controlled. For example, some patients may be more susceptible to a disease than others due to their dietary habits. Randomizing patients into the treatment or control group helps even out such differences, and it also prevents accidental bias from entering the study.

**Replication.** The more cases researchers observe, the more accurately they can estimate the effect of the explanatory variable on the response. In a single study, we **replicate** by collecting a sufficiently large sample. Additionally, a group of scientists may replicate an entire study to verify an earlier finding.

**Blocking.** Researchers sometimes know or suspect that variables, other than the treatment, influence the response. Under these circumstances, they may first group individuals based on this variable into **blocks** and then randomize cases within each block to the treatment groups. This strategy is often referred to as **blocking**. For instance, if we are looking at the effect of a drug on heart attacks, we might first split patients in the study into low-risk and high-risk blocks, then randomly assign half the patients from each block to the control group and the other half to the treatment group, as shown in Figure 1.16. This strategy ensures each treatment group has an equal number of low-risk and high-risk patients.

It is important to incorporate the first three experimental design principles into any study, and this book describes applicable methods for analyzing data from such experiments. Blocking is a slightly more advanced technique, and statistical methods in this book may be extended to analyze data collected using blocking.

### 1.5.2 Reducing bias in human experiments

Randomized experiments are the gold standard for data collection, but they do not ensure an unbiased perspective into the cause and effect relationships in all cases. Human studies are perfect examples where bias can unintentionally arise. Here we reconsider a study where a new drug was used to treat heart attack patients.<sup>16</sup> In particular, researchers wanted to know if the drug reduced deaths in patients.

---

<sup>16</sup>Anturane Reinfarction Trial Research Group. 1980. Sulfapyrazone in the prevention of sudden death after myocardial infarction. New England Journal of Medicine 302(5):250-256.

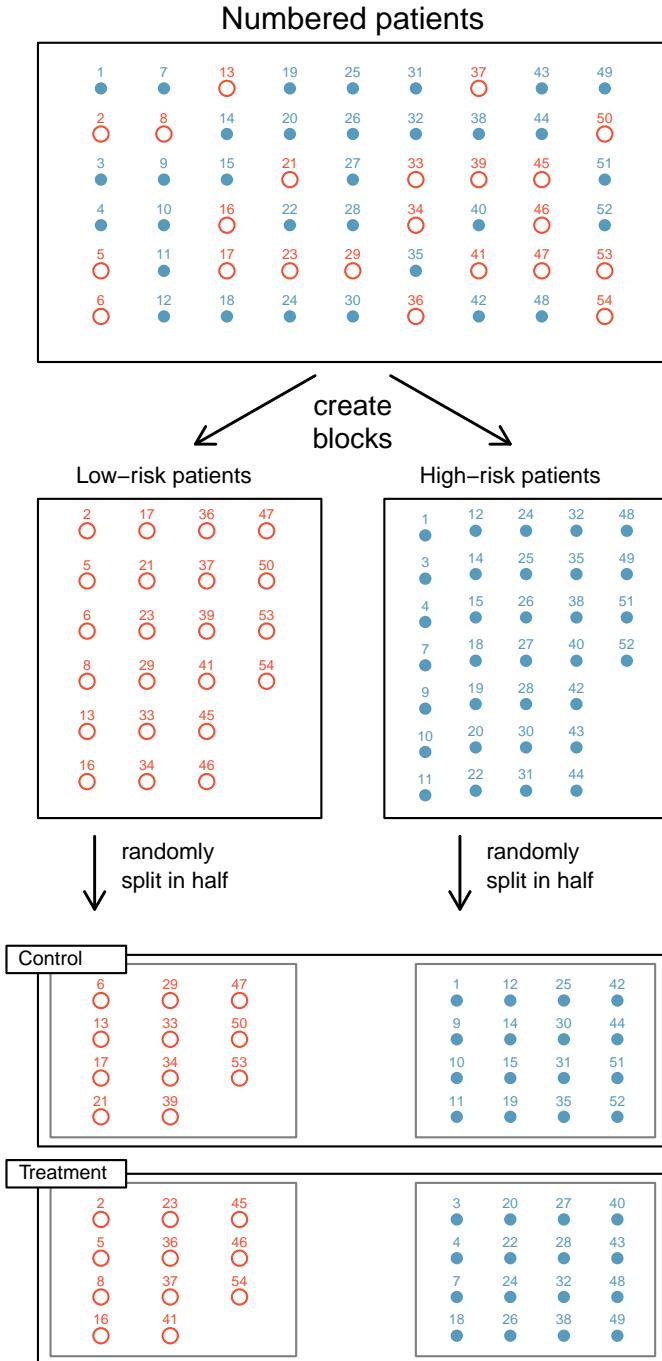


Figure 1.16: Blocking using a variable depicting patient risk. Patients are first divided into low-risk and high-risk blocks, then each block is evenly separated into the treatment groups using randomization. This strategy ensures an equal representation of patients in each treatment group from both the low-risk and high-risk categories.

These researchers designed a randomized experiment because they wanted to draw causal conclusions about the drug's effect. Study volunteers<sup>17</sup> were randomly placed into two study groups. One group, the **treatment group**, received the drug. The other group, called the **control group**, did not receive any drug treatment.

Put yourself in the place of a person in the study. If you are in the treatment group, you are given a fancy new drug that you anticipate will help you. On the other hand, a person in the other group doesn't receive the drug and sits idly, hoping her participation doesn't increase her risk of death. These perspectives suggest there are actually two effects: the one of interest is the effectiveness of the drug, and the second is an emotional effect that is difficult to quantify.

Researchers aren't usually interested in the emotional effect, which might bias the study. To circumvent this problem, researchers do not want patients to know which group they are in. When researchers keep the patients uninformed about their treatment, the study is said to be **blind**. But there is one problem: if a patient doesn't receive a treatment, she will know she is in the control group. The solution to this problem is to give fake treatments to patients in the control group. A fake treatment is called a **placebo**, and an effective placebo is the key to making a study truly blind. A classic example of a placebo is a sugar pill that is made to look like the actual treatment pill. Often times, a placebo results in a slight but real improvement in patients. This effect has been dubbed the **placebo effect**.

The patients are not the only ones who should be blinded: doctors and researchers can accidentally bias a study. When a doctor knows a patient has been given the real treatment, she might inadvertently give that patient more attention or care than a patient that she knows is on the placebo. To guard against this bias, which again has been found to have a measurable effect in some instances, most modern studies employ a **double-blind** setup where doctors or researchers who interact with patients are, just like the patients, unaware of who is or is not receiving the treatment.<sup>18</sup>

 **Guided Practice 1.14** Look back to the study in Section 1.1 where researchers were testing whether stents were effective at reducing strokes in at-risk patients. Is this an experiment? Was the study blinded? Was it double-blinded?<sup>19</sup>

## 1.6 Examining numerical data

In this section we will be introduced to techniques for exploring and summarizing numerical variables. The `email50` and `county` data sets from Section 1.2 provide rich opportunities for examples. Recall that outcomes of numerical variables are numbers on which it is reasonable to perform basic arithmetic operations. For example, the `pop2010` variable, which represents the populations of counties in 2010, is numerical since we can sensibly discuss the difference or ratio of the populations in two counties. On the other hand, area codes and zip codes are not numerical, but rather they are categorical variables.

---

<sup>17</sup>Human subjects are often called **patients**, **volunteers**, or **study participants**.

<sup>18</sup>There are always some researchers involved in the study who do know which patients are receiving which treatment. However, they do not interact with the study's patients and do not tell the blinded health care professionals who is receiving which treatment.

<sup>19</sup>The researchers assigned the patients into their treatment groups, so this study was an experiment. However, the patients could distinguish what treatment they received, so this study was not blind. The study could not be double-blind since it was not blind.

### 1.6.1 Scatterplots for paired data

A **scatterplot** provides a case-by-case view of data for two numerical variables. In Figure 1.8 on page 13, a scatterplot was used to examine how federal spending and poverty were related in the `county` data set. Another scatterplot is shown in Figure 1.17, comparing the number of line breaks (`line_breaks`) and number of characters (`num_char`) in emails for the `email150` data set. In any scatterplot, each point represents a single case. Since there are 50 cases in `email150`, there are 50 points in Figure 1.17.

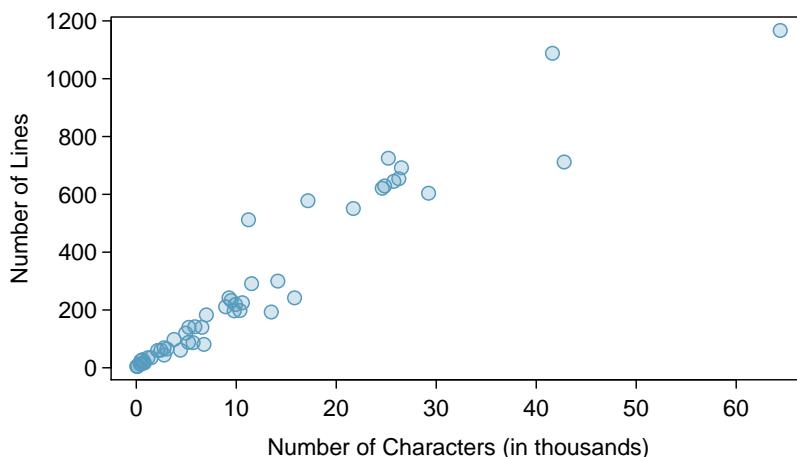


Figure 1.17: A scatterplot of `line_breaks` versus `num_char` for the `email150` data.

To put the number of characters in perspective, this paragraph has 363 characters. Looking at Figure 1.17, it seems that some emails are incredibly verbose! Upon further investigation, we would actually find that most of the long emails use the HTML format, which means most of the characters in those emails are used to format the email rather than provide text.

Ⓐ **Guided Practice 1.15** What do scatterplots reveal about the data, and how might they be useful?<sup>20</sup>

Ⓑ **Example 1.16** Consider a new data set of 54 cars with two variables: vehicle price and weight.<sup>21</sup> A scatterplot of vehicle price versus weight is shown in Figure 1.18. What can be said about the relationship between these variables?

---

The relationship is evidently nonlinear, as highlighted by the dashed line. This is different from previous scatterplots we've seen, such as Figure 1.8 on page 13 and Figure 1.17, which show relationships that are very linear.

Ⓐ **Guided Practice 1.17** Describe two variables that would have a horseshoe shaped association in a scatterplot.<sup>22</sup>

<sup>20</sup> Answers may vary. Scatterplots are helpful in quickly spotting associations relating variables, whether those associations come in the form of simple trends or whether those relationships are more complex.

<sup>21</sup>Subset of data from [www.amstat.org/publications/jse/v1n1/datasets.lock.html](http://www.amstat.org/publications/jse/v1n1/datasets.lock.html)

<sup>22</sup>Consider the case where your vertical axis represents something “good” and your horizontal axis represents something that is only good in moderation. Health and water consumption fit this description since water becomes toxic when consumed in excessive quantities.

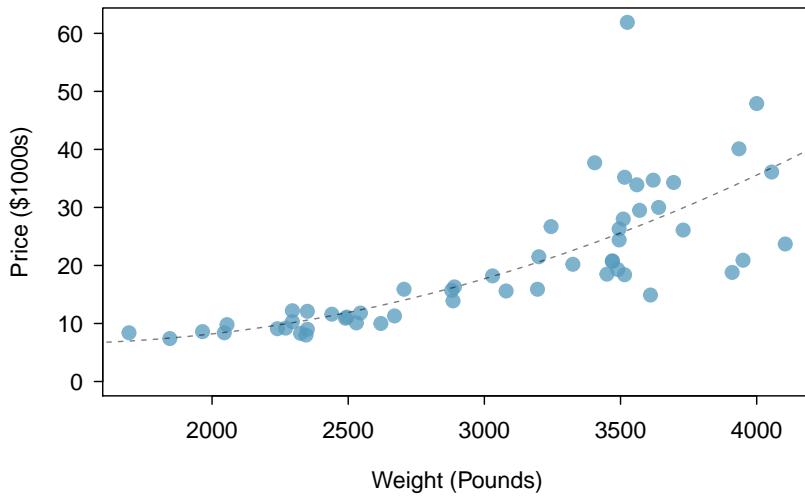


Figure 1.18: A scatterplot of `price` versus `weight` for 54 cars.

### 1.6.2 Dot plots and the mean

Sometimes two variables are one too many: only one variable may be of interest. In these cases, a dot plot provides the most basic of displays. A **dot plot** is a one-variable scatterplot; an example using the number of characters from 50 emails is shown in Figure 1.19. A stacked version of this dot plot is shown in Figure 1.20.

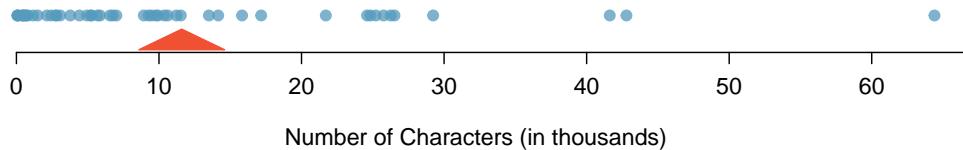


Figure 1.19: A dot plot of `num_char` for the `email150` data set.

The **mean**, sometimes called the average, is a common way to measure the center of a **distribution** of data. To find the mean number of characters in the 50 emails, we add up all the character counts and divide by the number of emails. For computational convenience, the number of characters is listed in the thousands and rounded to the first decimal.

$$\bar{x} = \frac{21.7 + 7.0 + \dots + 15.8}{50} = 11.6 \quad (1.18)$$

$\bar{x}$   
sample  
mean

The sample mean is often labeled  $\bar{x}$ . The letter  $x$  is being used as a generic placeholder for the variable of interest, `num_char`, and the bar over on the  $x$  communicates that the average number of characters in the 50 emails was 11,600. It is useful to think of the mean as the balancing point of the distribution. The sample mean is shown as a triangle in Figures 1.19 and 1.20.

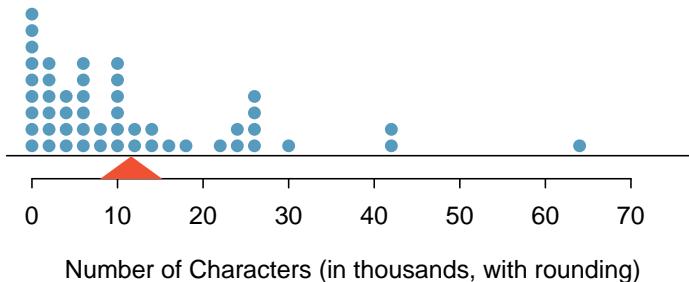


Figure 1.20: A stacked dot plot of `num_char` for the `email150` data set.  
The values have been rounded to the nearest 2,000 in this plot.

### Mean

The sample mean of a numerical variable is computed as the sum of all of the observations divided by the number of observations:

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n} \quad (1.19)$$

where  $x_1, x_2, \dots, x_n$  represent the  $n$  observed values.

$n$   
sample size

- Ⓐ **Guided Practice 1.20** Examine Equations (1.18) and (1.19) above. What does  $x_1$  correspond to? And  $x_2$ ? Can you infer a general meaning to what  $x_i$  might represent?<sup>23</sup>

- Ⓑ **Guided Practice 1.21** What was  $n$  in this sample of emails?<sup>24</sup>

The `email150` data set represents a sample from a larger population of emails that were received in January and March. We could compute a mean for this population in the same way as the sample mean. However, the population mean has a special label:  $\mu$ . The symbol  $\mu$  is the Greek letter *mu* and represents the average of all observations in the population. Sometimes a subscript, such as  $_x$ , is used to represent which variable the population mean refers to, e.g.  $\mu_x$ .

$\mu$   
population  
mean

- Ⓒ **Example 1.22** The average number of characters across all emails can be estimated using the sample data. Based on the sample of 50 emails, what would be a reasonable estimate of  $\mu_x$ , the mean number of characters in all emails in the `email` data set? (Recall that `email150` is a sample from `email`.)

The sample mean, 11,600, may provide a reasonable estimate of  $\mu_x$ . While this number will not be perfect, it provides a *point estimate* of the population mean. In Chapter 4 and beyond, we will develop tools to characterize the accuracy of point estimates, and we will find that point estimates based on larger samples tend to be more accurate than those based on smaller samples.

<sup>23</sup> $x_1$  corresponds to the number of characters in the first email in the sample (21.7, in thousands),  $x_2$  to the number of characters in the second email (7.0, in thousands), and  $x_i$  corresponds to the number of characters in the  $i^{th}$  email in the data set.

<sup>24</sup>The sample size was  $n = 50$ .

- **Example 1.23** We might like to compute the average income per person in the US. To do so, we might first think to take the mean of the per capita incomes across the 3,143 counties in the `county` data set. What would be a better approach?

The `county` data set is special in that each county actually represents many individual people. If we were to simply average across the `income` variable, we would be treating counties with 5,000 and 5,000,000 residents equally in the calculations. Instead, we should compute the total income for each county, add up all the counties' totals, and then divide by the number of people in all the counties. If we completed these steps with the `county` data, we would find that the per capita income for the US is \$27,348.43. Had we computed the *simple* mean of per capita income across counties, the result would have been just \$22,504.70!

Example 1.23 used what is called a **weighted mean**, which will not be a key topic in this textbook. However, we have provided an online supplement on weighted means for interested readers:

[www.openintro.org/d?f=wtdmean](http://www.openintro.org/d?f=wtdmean)

### 1.6.3 Histograms and shape

Dot plots show the exact value for each observation. This is useful for small data sets, but they can become hard to read with larger samples. Rather than showing the value of each observation, we prefer to think of the value as belonging to a *bin*. For example, in the `email150` data set, we create a table of counts for the number of cases with character counts between 0 and 5,000, then the number of cases between 5,000 and 10,000, and so on. Observations that fall on the boundary of a bin (e.g. 5,000) are allocated to the lower bin. This tabulation is shown in Table 1.21. These binned counts are plotted as bars in Figure 1.22 into what is called a **histogram**, which resembles the stacked dot plot shown in Figure 1.20.

Characters (in thousands)	0-5	5-10	10-15	15-20	20-25	25-30	...	55-60	60-65
Count	19	12	6	2	3	5	...	0	1

Table 1.21: The counts for the binned `num_char` data.

Histograms provide a view of the **data density**. Higher bars represent where the data are relatively more common. For instance, there are many more emails with fewer than 20,000 characters than emails with at least 20,000 in the data set. The bars make it easy to see how the density of the data changes relative to the number of characters.

Histograms are especially convenient for describing the shape of the data distribution. Figure 1.22 shows that most emails have a relatively small number of characters, while fewer emails have a very large number of characters. When data trail off to the right in this way and have a longer right tail, the shape is said to be **right skewed**.<sup>25</sup>

Data sets with the reverse characteristic – a long, thin tail to the left – are said to be **left skewed**. We also say that such a distribution has a long left tail. Data sets that show roughly equal trailing off in both directions are called **symmetric**.

<sup>25</sup>Other ways to describe data that are skewed to the right: **skewed to the right**, **skewed to the high end**, or **skewed to the positive end**.

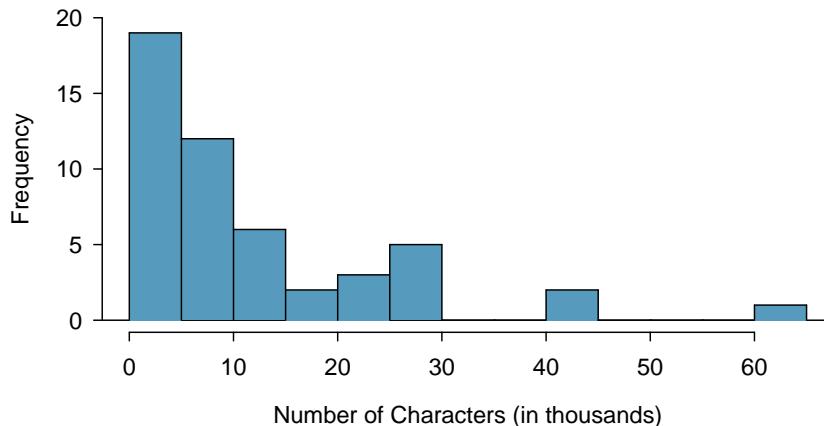


Figure 1.22: A histogram of `num_char`. This distribution is very strongly skewed to the right.

#### Long tails to identify skew

When data trail off in one direction, the distribution has a **long tail**. If a distribution has a long left tail, it is left skewed. If a distribution has a long right tail, it is right skewed.

Ⓐ **Guided Practice 1.24** Take a look at the dot plots in Figures 1.19 and 1.20. Can you see the skew in the data? Is it easier to see the skew in this histogram or the dot plots?<sup>26</sup>

Ⓑ **Guided Practice 1.25** Besides the mean (since it was labeled), what can you see in the dot plots that you cannot see in the histogram?<sup>27</sup>

In addition to looking at whether a distribution is skewed or symmetric, histograms can be used to identify modes. A **mode** is represented by a prominent peak in the distribution.<sup>28</sup> There is only one prominent peak in the histogram of `num_char`.

Figure 1.23 shows histograms that have one, two, or three prominent peaks. Such distributions are called **unimodal**, **bimodal**, and **multimodal**, respectively. Any distribution with more than 2 prominent peaks is called multimodal. Notice that there was one prominent peak in the unimodal distribution with a second less prominent peak that was not counted since it only differs from its neighboring bins by a few observations.

Ⓒ **Guided Practice 1.26** Figure 1.22 reveals only one prominent mode in the number of characters. Is the distribution unimodal, bimodal, or multimodal?<sup>29</sup>

<sup>26</sup>The skew is visible in all three plots, though the flat dot plot is the least useful. The stacked dot plot and histogram are helpful visualizations for identifying skew.

<sup>27</sup>Character counts for individual emails.

<sup>28</sup>Another definition of mode, which is not typically used in statistics, is the value with the most occurrences. It is common to have *no* observations with the same value in a data set, which makes this other definition useless for many real data sets.

<sup>29</sup>Unimodal. Remember that *uni* stands for 1 (think *unicycles*). Similarly, *bi* stands for 2 (think *bicycles*). (We're hoping a *multicycle* will be invented to complete this analogy.)

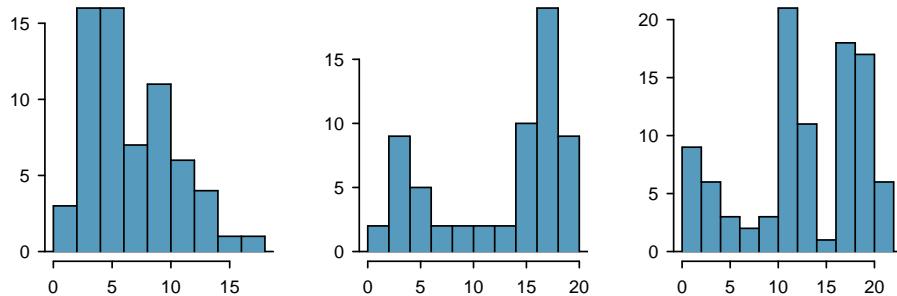


Figure 1.23: Counting only prominent peaks, the distributions are (left to right) unimodal, bimodal, and multimodal. Note that we've said the left plot is unimodal intentionally. This is because we are counting *prominent* peaks, not just any peak. See the Tip box for more thoughts.

- Ⓐ **Guided Practice 1.27** Height measurements of young students and adult teachers at a K-3 elementary school were taken. How many modes would you anticipate in this height data set?<sup>30</sup>

**TIP: Looking for modes**

Looking for modes isn't about finding a clear and correct answer about the number of modes in a distribution, which is why *prominent* is not rigorously defined in this book. The important part of this examination is to better understand your data and how it might be structured.

#### 1.6.4 Variance and standard deviation

The mean was introduced as a method to describe the center of a data set, but the variability in the data is also important. Here, we introduce two measures of variability: the variance and the standard deviation. Both of these are very useful in data analysis, even though their formulas are a bit tedious to calculate by hand. The standard deviation is the easier of the two to understand, and it roughly describes how far away the typical observation is from the mean.

We call the distance of an observation from its mean its **deviation**. Below are the deviations for the 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup>, and 50<sup>th</sup> observations in the `num_char` variable. For computational convenience, the number of characters is listed in the thousands and rounded to the first decimal.

$$x_1 - \bar{x} = 21.7 - 11.6 = 10.1$$

$$x_2 - \bar{x} = 7.0 - 11.6 = -4.6$$

$$x_3 - \bar{x} = 0.6 - 11.6 = -11.0$$

⋮

$$x_{50} - \bar{x} = 15.8 - 11.6 = 4.2$$

<sup>30</sup>There might be two height groups visible in the data set: one of the students and one of the adults. That is, the data are probably bimodal.

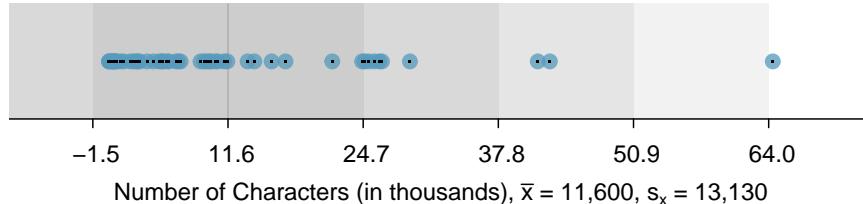


Figure 1.24: In the `num_char` data, 41 of the 50 emails (82%) are within 1 standard deviation of the mean, and 47 of the 50 emails (94%) are within 2 standard deviations. Usually about 70% of the data are within 1 standard deviation of the mean and 95% are within 2 standard deviations, though this rule of thumb is less accurate for skewed data, as shown in this example.

If we square these deviations and then take an average, the result is about equal to the sample **variance**, denoted by  $s^2$ :

$$\begin{aligned}s^2 &= \frac{10.1^2 + (-4.6)^2 + (-11.0)^2 + \dots + 4.2^2}{50 - 1} \\ &= \frac{102.01 + 21.16 + 121.00 + \dots + 17.64}{49} \\ &= 172.44\end{aligned}$$

$s^2$   
sample variance

We divide by  $n - 1$ , rather than dividing by  $n$ , when computing the variance; you need not worry about this mathematical nuance for the material in this textbook. Notice that squaring the deviations does two things. First, it makes large values much larger, seen by comparing  $10.1^2$ ,  $(-4.6)^2$ ,  $(-11.0)^2$ , and  $4.2^2$ . Second, it gets rid of any negative signs.

The **standard deviation** is defined as the square root of the variance:

$$s = \sqrt{172.44} = 13.13$$

$s$   
sample standard deviation

The standard deviation of the number of characters in an email is about 13.13 thousand. A subscript of  $_x$  may be added to the variance and standard deviation, i.e.  $s_x^2$  and  $s_x$ , as a reminder that these are the variance and standard deviation of the observations represented by  $x_1, x_2, \dots, x_n$ . The  $_x$  subscript is usually omitted when it is clear which data the variance or standard deviation is referencing.

### Variance and standard deviation

The variance is roughly the average squared distance from the mean. The standard deviation is the square root of the variance. The standard deviation is useful when considering how close the data are to the mean.

Formulas and methods used to compute the variance and standard deviation for a population are similar to those used for a sample.<sup>31</sup> However, like the mean, the population values have special symbols:  $\sigma^2$  for the variance and  $\sigma$  for the standard deviation. The symbol  $\sigma$  is the Greek letter *sigma*.

$\sigma^2$   
population variance  
 $\sigma$   
population standard deviation

<sup>31</sup>The only difference is that the population variance has a division by  $n$  instead of  $n - 1$ .

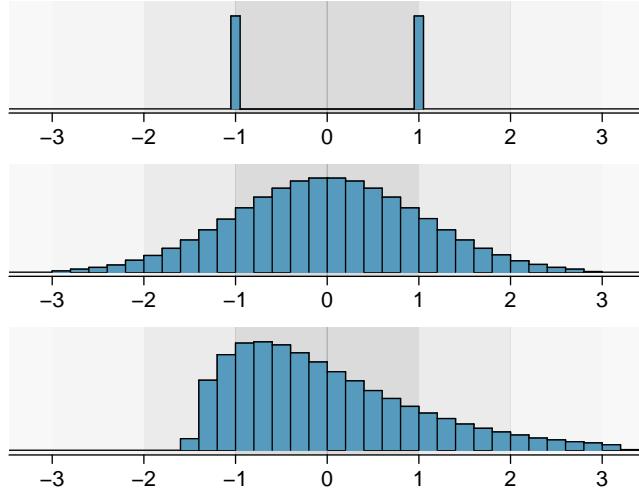


Figure 1.25: Three very different population distributions with the same mean  $\mu = 0$  and standard deviation  $\sigma = 1$ .

**TIP: standard deviation describes variability**

Focus on the conceptual meaning of the standard deviation as a descriptor of variability rather than the formulas. Usually 70% of the data will be within one standard deviation of the mean and about 95% will be within two standard deviations. However, as seen in Figures 1.24 and 1.25, these percentages are not strict rules.

- Ⓐ **Guided Practice 1.28** On page 30, the concept of shape of a distribution was introduced. A good description of the shape of a distribution should include modality and whether the distribution is symmetric or skewed to one side. Using Figure 1.25 as an example, explain why such a description is important.<sup>32</sup>

- Ⓑ **Example 1.29** Describe the distribution of the `num_char` variable using the histogram in Figure 1.22 on page 31. The description should incorporate the center, variability, and shape of the distribution, and it should also be placed in context: the number of characters in emails. Also note any especially unusual cases.

The distribution of email character counts is unimodal and very strongly skewed to the high end. Many of the counts fall near the mean at 11,600, and most fall within one standard deviation (13,130) of the mean. There is one exceptionally long email with about 65,000 characters.

In practice, the variance and standard deviation are sometimes used as a means to an end, where the “end” is being able to accurately estimate the uncertainty associated with a sample statistic. For example, in Chapter 4 we will use the variance and standard deviation to assess how close the sample mean is to the population mean.

<sup>32</sup>Figure 1.25 shows three distributions that look quite different, but all have the same mean, variance, and standard deviation. Using modality, we can distinguish between the first plot (bimodal) and the last two (unimodal). Using skewness, we can distinguish between the last plot (right skewed) and the first two. While a picture, like a histogram, tells a more complete story, we can use modality and shape (symmetry/skew) to characterize basic information about a distribution.

### 1.6.5 Box plots, quartiles, and the median

A **box plot** summarizes a data set using five statistics while also plotting unusual observations. Figure 1.26 provides a vertical dot plot alongside a box plot of the `num_char` variable from the `email150` data set.

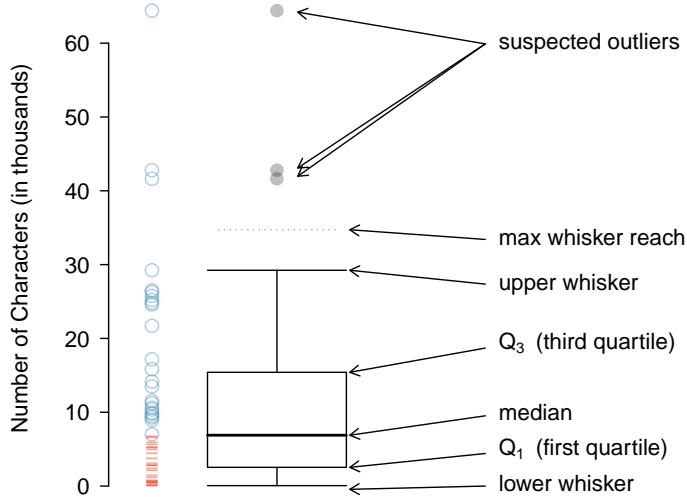


Figure 1.26: A vertical dot plot next to a labeled box plot for the number of characters in 50 emails. The median (6,890), splits the data into the bottom 50% and the top 50%, marked in the dot plot by horizontal dashes and open circles, respectively.

The first step in building a box plot is drawing a dark line denoting the **median**, which splits the data in half. Figure 1.26 shows 50% of the data falling below the median (dashes) and other 50% falling above the median (open circles). There are 50 character counts in the data set (an even number) so the data are perfectly split into two groups of 25. We take the median in this case to be the average of the two observations closest to the 50<sup>th</sup> percentile:  $(6,768 + 7,012)/2 = 6,890$ . When there are an odd number of observations, there will be exactly one observation that splits the data into two halves, and in such a case that observation is the median (no average needed).

#### Median: the number in the middle

If the data are ordered from smallest to largest, the **median** is the observation right in the middle. If there are an even number of observations, there will be two values in the middle, and the median is taken as their average.

The second step in building a box plot is drawing a rectangle to represent the middle 50% of the data. The total length of the box, shown vertically in Figure 1.26, is called the **interquartile range** (IQR, for short). It, like the standard deviation, is a measure of variability in data. The more variable the data, the larger the standard deviation and IQR. The two boundaries of the box are called the **first quartile** (the 25<sup>th</sup> percentile, i.e. 25% of the data fall below this value) and the **third quartile** (the 75<sup>th</sup> percentile), and these are often labeled  $Q_1$  and  $Q_3$ , respectively.

### Interquartile range (IQR)

The IQR is the length of the box in a box plot. It is computed as

$$IQR = Q_3 - Q_1$$

where  $Q_1$  and  $Q_3$  are the 25<sup>th</sup> and 75<sup>th</sup> percentiles.

- **Guided Practice 1.30** What percent of the data fall between  $Q_1$  and the median? What percent is between the median and  $Q_3$ ?<sup>33</sup>

Extending out from the box, the **whiskers** attempt to capture the data outside of the box, however, their reach is never allowed to be more than  $1.5 \times IQR$ .<sup>34</sup> They capture everything within this reach. In Figure 1.26, the upper whisker does not extend to the last three points, which is beyond  $Q_3 + 1.5 \times IQR$ , and so it extends only to the last point below this limit. The lower whisker stops at the lowest value, 33, since there is no additional data to reach; the lower whisker's limit is not shown in the figure because the plot does not extend down to  $Q_1 - 1.5 \times IQR$ . In a sense, the box is like the body of the box plot and the whiskers are like its arms trying to reach the rest of the data.

Any observation that lies beyond the whiskers is labeled with a dot. The purpose of labeling these points – instead of just extending the whiskers to the minimum and maximum observed values – is to help identify any observations that appear to be unusually distant from the rest of the data. Unusually distant observations are called **outliers**. In this case, it would be reasonable to classify the emails with character counts of 41,623, 42,793, and 64,401 as outliers since they are numerically distant from most of the data.

### Outliers are extreme

An **outlier** is an observation that appears extreme relative to the rest of the data.

### TIP: Why it is important to look for outliers

Examination of data for possible outliers serves many useful purposes, including

1. Identifying strong skew in the distribution.
2. Identifying data collection or entry errors. For instance, we re-examined the email purported to have 64,401 characters to ensure this value was accurate.
3. Providing insight into interesting properties of the data.

- **Guided Practice 1.31** The observation 64,401, a suspected outlier, was found to be an accurate observation. What would such an observation suggest about the nature of character counts in emails?<sup>35</sup>

<sup>33</sup>Since  $Q_1$  and  $Q_3$  capture the middle 50% of the data and the median splits the data in the middle, 25% of the data fall between  $Q_1$  and the median, and another 25% falls between the median and  $Q_3$ .

<sup>34</sup>While the choice of exactly 1.5 is arbitrary, it is the most commonly used value for box plots.

<sup>35</sup>That occasionally there may be very long emails.

- Ⓐ **Guided Practice 1.32** Using Figure 1.26, estimate the following values for `num_char` in the `email150` data set: (a)  $Q_1$ , (b)  $Q_3$ , and (c) IQR.<sup>36</sup>

 **Calculator videos**

Videos covering how to create statistical summaries and box plots using TI and Casio graphing calculators are available at [openintro.org/videos](http://openintro.org/videos).

### 1.6.6 Robust statistics

How are the sample statistics of the `num_char` data set affected by the observation, 64,401? What would have happened if this email wasn't observed? What would happen to these summary statistics if the observation at 64,401 had been even larger, say 150,000? These scenarios are plotted alongside the original data in Figure 1.27, and sample statistics are computed under each scenario in Table 1.28.

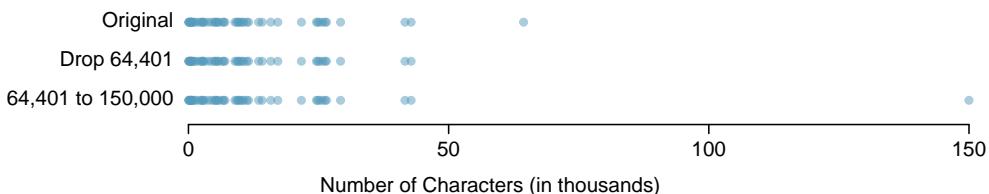


Figure 1.27: Dot plots of the original character count data and two modified data sets.

scenario	robust		not robust	
	median	IQR	$\bar{x}$	$s$
original <code>num_char</code> data	6,890	12,875	11,600	13,130
drop 64,401 observation	6,768	11,702	10,521	10,798
move 64,401 to 150,000	6,890	12,875	13,310	22,434

Table 1.28: A comparison of how the median, IQR, mean ( $\bar{x}$ ), and standard deviation ( $s$ ) change when extreme observations are present.

- Ⓐ **Guided Practice 1.33** (a) Which is more affected by extreme observations, the mean or median? Table 1.28 may be helpful. (b) Is the standard deviation or IQR more affected by extreme observations?<sup>37</sup>

The median and IQR are called **robust estimates** because extreme observations have little effect on their values. The mean and standard deviation are much more affected by changes in extreme observations.

<sup>36</sup>These visual estimates will vary a little from one person to the next:  $Q_1 = 3,000$ ,  $Q_3 = 15,000$ ,  $IQR = Q_3 - Q_1 = 12,000$ . (The true values:  $Q_1 = 2,536$ ,  $Q_3 = 15,411$ ,  $IQR = 12,875$ .)

<sup>37</sup>(a) Mean is affected more. (b) Standard deviation is affected more. Complete explanations are provided in the material following Guided Practice 1.33.

- **Example 1.34** The median and IQR do not change much under the three scenarios in Table 1.28. Why might this be the case?

The median and IQR are only sensitive to numbers near  $Q_1$ , the median, and  $Q_3$ . Since values in these regions are relatively stable – there aren’t large jumps between observations – the median and IQR estimates are also quite stable.

- **Guided Practice 1.35** The distribution of vehicle prices tends to be right skewed, with a few luxury and sports cars lingering out into the right tail. If you were searching for a new car and cared about price, should you be more interested in the mean or median price of vehicles sold, assuming you are in the market for a regular car?<sup>38</sup>

### 1.6.7 Transforming data (special topic)

When data are very strongly skewed, we sometimes transform them so they are easier to model. Consider the histogram of salaries for Major League Baseball players’ salaries from 2010, which is shown in Figure 1.29(a).

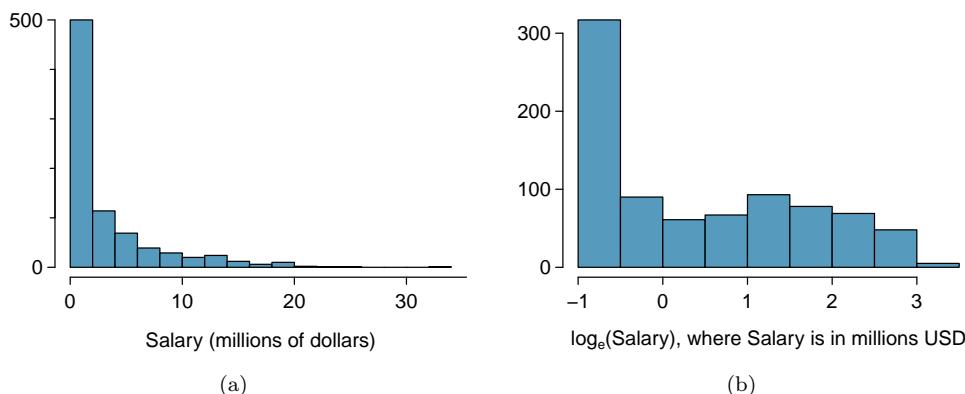


Figure 1.29: (a) Histogram of MLB player salaries for 2010, in millions of dollars. (b) Histogram of the log-transformed MLB player salaries for 2010.

- **Example 1.36** The histogram of MLB player salaries is useful in that we can see the data are extremely skewed and centered (as gauged by the median) at about \$1 million. What isn’t useful about this plot?

Most of the data are collected into one bin in the histogram and the data are so strongly skewed that many details in the data are obscured.

There are some standard transformations that are often applied when much of the data cluster near zero (relative to the larger values in the data set) and all observations are positive. A **transformation** is a rescaling of the data using a function. For instance, a plot of the natural logarithm<sup>39</sup> of player salaries results in a new histogram in Figure 1.29(b).

<sup>38</sup>Buyers of a “regular car” should be concerned about the median price. High-end car sales can drastically inflate the mean price while the median will be more robust to the influence of those sales.

<sup>39</sup>Statisticians often write the natural logarithm as  $\log$ . You might be more familiar with it being written as  $\ln$ .

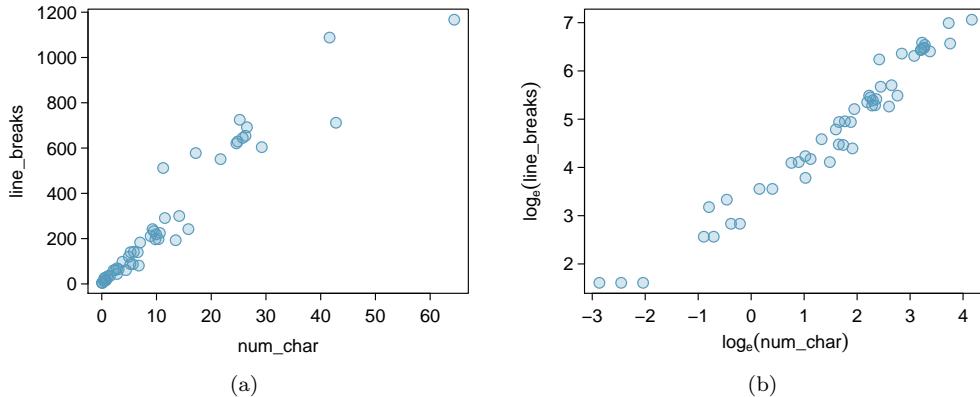


Figure 1.30: (a) Scatterplot of `line_breaks` against `num_char` for 50 emails.  
 (b) A scatterplot of the same data but where each variable has been log-transformed.

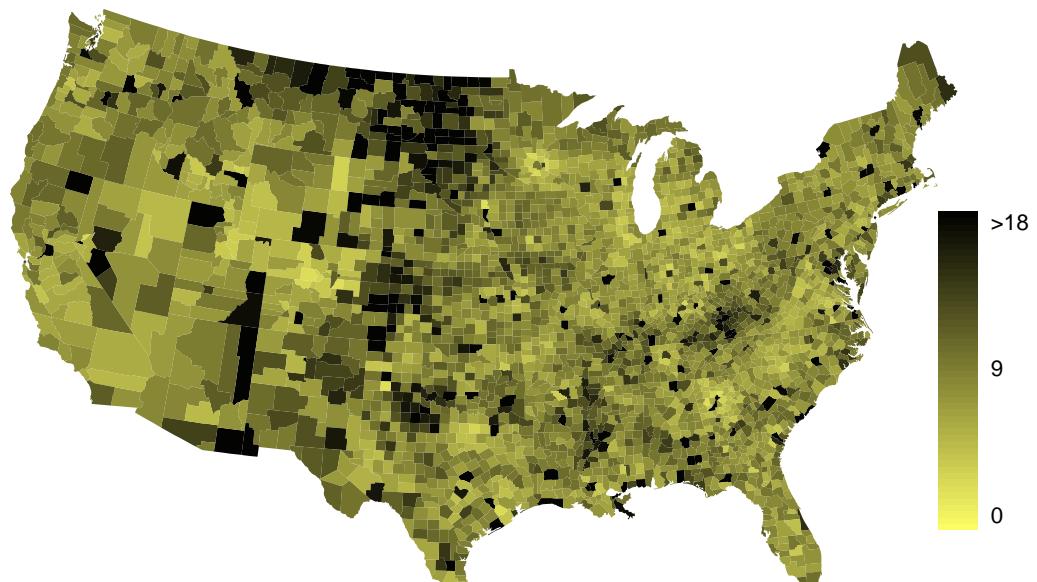
Transformed data are sometimes easier to work with when applying statistical models because the transformed data are much less skewed and outliers are usually less extreme.

Transformations can also be applied to one or both variables in a scatterplot. A scatterplot of the `line_breaks` and `num_char` variables is shown in Figure 1.30(a), which was earlier shown in Figure 1.17. We can see a positive association between the variables and that many observations are clustered near zero. In Chapter 7, we might want to use a straight line to model the data. However, we'll find that the data in their current state cannot be modeled very well. Figure 1.30(b) shows a scatterplot where both the `line_breaks` and `num_char` variables have been transformed using a log (base  $e$ ) transformation. While there is a positive association in each plot, the transformed data show a steadier trend, which is easier to model than the untransformed data.

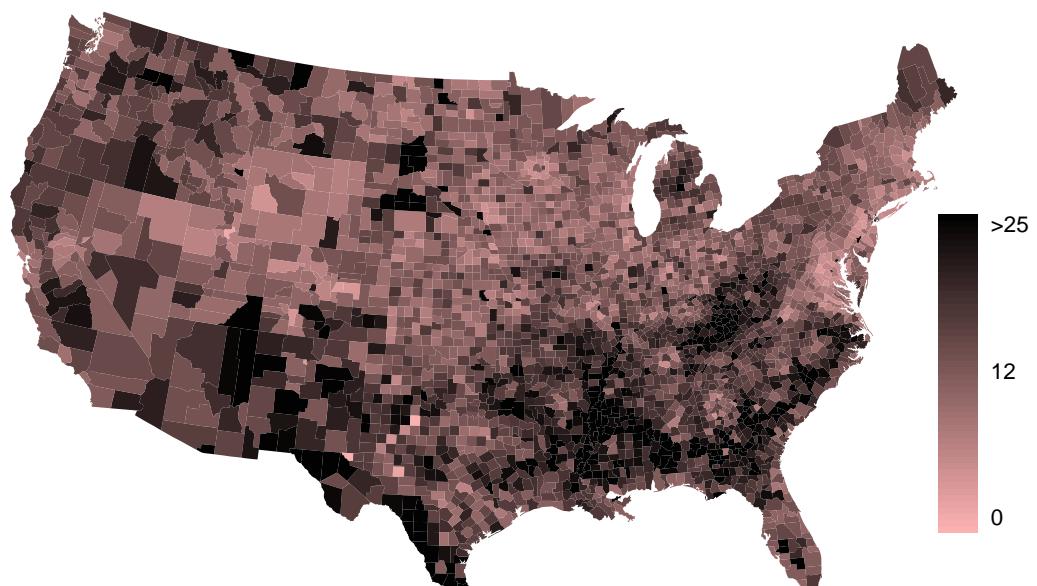
Transformations other than the logarithm can be useful, too. For instance, the square root ( $\sqrt{\text{original observation}}$ ) and inverse ( $\frac{1}{\text{original observation}}$ ) are used by statisticians. Common goals in transforming data are to see the data structure differently, reduce skew, assist in modeling, or straighten a nonlinear relationship in a scatterplot.

## 1.6.8 Mapping data (special topic)

The `county` data set offers many numerical variables that we could plot using dot plots, scatterplots, or box plots, but these miss the true nature of the data. Rather, when we encounter geographic data, we should map it using an **intensity map**, where colors are used to show higher and lower values of a variable. Figures 1.31 and 1.32 shows intensity maps for federal spending per capita (`fed_spend`), poverty rate in percent (`poverty`), homeownership rate in percent (`homeownership`), and median household income (`med_income`). The color key indicates which colors correspond to which values. Note that the intensity maps are not generally very helpful for getting precise values in any given county, but they are very helpful for seeing geographic trends and generating interesting research questions.



(a)



(b)

Figure 1.31: (a) Map of federal spending (dollars per capita). (b) Intensity map of poverty rate (percent).

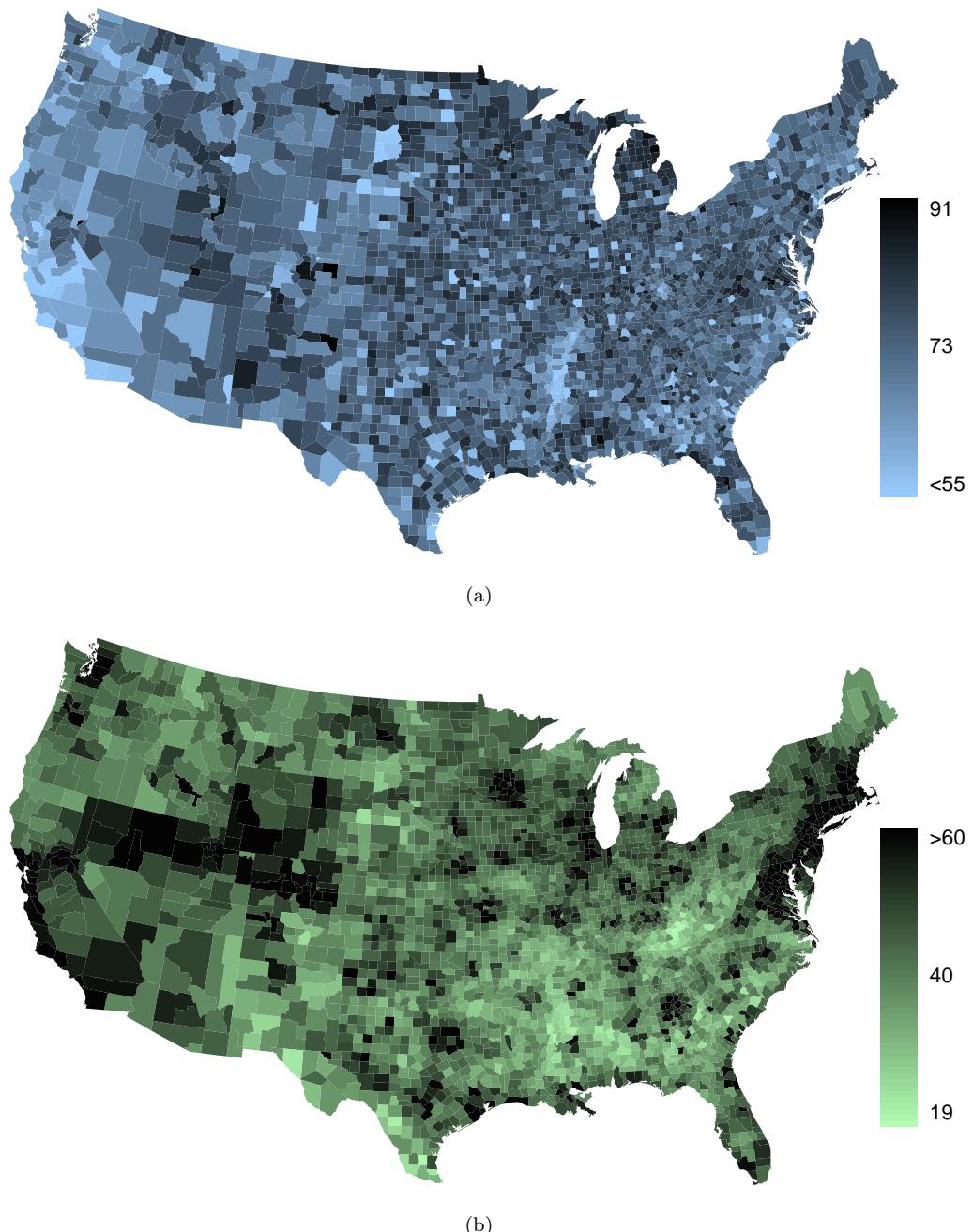


Figure 1.32: (a) Intensity map of homeownership rate (percent). (b) Intensity map of median household income (\$1000s).

- **Example 1.37** What interesting features are evident in the `fed_spend` and `poverty` intensity maps?

The federal spending intensity map shows substantial spending in the Dakotas and along the central-to-western part of the Canadian border, which may be related to the oil boom in this region. There are several other patches of federal spending, such as a vertical strip in eastern Utah and Arizona and the area where Colorado, Nebraska, and Kansas meet. There are also seemingly random counties with very high federal spending relative to their neighbors. If we did not cap the federal spending range at \$18 per capita, we would actually find that some counties have extremely high federal spending while there is almost no federal spending in the neighboring counties. These high-spending counties might contain military bases, companies with large government contracts, or other government facilities with many employees.

Poverty rates are evidently higher in a few locations. Notably, the deep south shows higher poverty rates, as does the southwest border of Texas. The vertical strip of eastern Utah and Arizona, noted above for its higher federal spending, also appears to have higher rates of poverty (though generally little correspondence is seen between the two variables). High poverty rates are evident in the Mississippi flood plains a little north of New Orleans and also in a large section of Kentucky and West Virginia.

- **Guided Practice 1.38** What interesting features are evident in the `med_income` intensity map in Figure 1.32(b)?<sup>40</sup>

---

<sup>40</sup>Note: answers will vary. There is a very strong correspondence between high earning and metropolitan areas. You might look for large cities you are familiar with and try to spot them on the map as dark spots.

## 1.7 Considering categorical data

Like numerical data, categorical data can also be organized and analyzed. In this section, we will introduce tables and other basic tools for categorical data that are used throughout this book. The `email150` data set represents a sample from a larger email data set called `email`. This larger data set contains information on 3,921 emails. In this section we will examine whether the presence of numbers, small or large, in an email provides any useful value in classifying email as spam or not spam.

### 1.7.1 Contingency tables and bar plots

Table 1.33 summarizes two variables: `spam` and `number`. Recall that `number` is a categorical variable that describes whether an email contains no numbers, only small numbers (values under 1 million), or at least one big number (a value of 1 million or more). A table that summarizes data for two categorical variables in this way is called a **contingency table**. Each value in the table represents the number of times a particular combination of variable outcomes occurred. For example, the value 149 corresponds to the number of emails in the data set that are spam *and* had no number listed in the email. Row and column totals are also included. The **row totals** provide the total counts across each row (e.g.  $149 + 168 + 50 = 367$ ), and **column totals** are total counts down each column.

A table for a single variable is called a **frequency table**. Table 1.34 is a frequency table for the `number` variable. If we replaced the counts with percentages or proportions, the table would be called a **relative frequency table**.

		number			
		none	small	big	Total
spam	spam	149	168	50	367
	not spam	400	2659	495	3554
	Total	549	2827	545	3921

Table 1.33: A contingency table for `spam` and `number`.

	none	small	big	Total
	549	2827	545	3921

Table 1.34: A frequency table for the `number` variable.

A bar plot is a common way to display a single categorical variable. The left panel of Figure 1.35 shows a **bar plot** for the `number` variable. In the right panel, the counts are converted into proportions (e.g.  $549/3921 = 0.140$  for `none`), showing the proportion of observations that are in each level (i.e. in each category).

### 1.7.2 Row and column proportions

Table 1.36 shows the row proportions for Table 1.33. The **row proportions** are computed as the counts divided by their row totals. The value 149 at the intersection of `spam` and `none` is replaced by  $149/367 = 0.406$ , i.e. 149 divided by its row total, 367. So what does 0.406 represent? It corresponds to the proportion of spam emails in the sample that do not have any numbers.

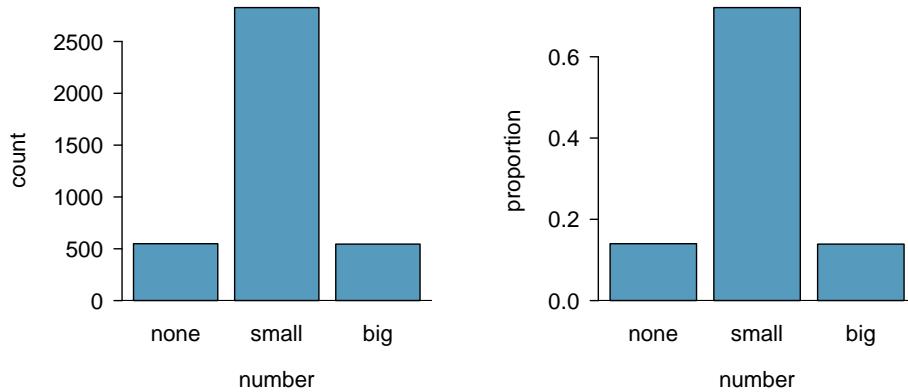


Figure 1.35: Two bar plots of `number`. The left panel shows the counts, and the right panel shows the proportions in each group.

	none	small	big	Total
spam	$149/367 = 0.406$	$168/367 = 0.458$	$50/367 = 0.136$	1.000
not spam	$400/3554 = 0.113$	$2659/3554 = 0.748$	$495/3554 = 0.139$	1.000
Total	$549/3921 = 0.140$	$2827/3921 = 0.721$	$545/3921 = 0.139$	1.000

Table 1.36: A contingency table with row proportions for the `spam` and `number` variables.

A contingency table of the column proportions is computed in a similar way, where each **column proportion** is computed as the count divided by the corresponding column total. Table 1.37 shows such a table, and here the value 0.271 indicates that 27.1% of emails with no numbers were spam. This rate of spam is much higher compared to emails with only small numbers (5.9%) or big numbers (9.2%). Because these spam rates vary between the three levels of `number` (`none`, `small`, `big`), this provides evidence that the `spam` and `number` variables are associated.

	none	small	big	Total
spam	$149/549 = 0.271$	$168/2827 = 0.059$	$50/545 = 0.092$	$367/3921 = 0.094$
not spam	$400/549 = 0.729$	$2659/2827 = 0.941$	$495/545 = 0.908$	$3554/3921 = 0.906$
Total	1.000	1.000	1.000	1.000

Table 1.37: A contingency table with column proportions for the `spam` and `number` variables.

We could also have checked for an association between `spam` and `number` in Table 1.36 using row proportions. When comparing these row proportions, we would look down columns to see if the fraction of emails with no numbers, small numbers, and big numbers varied from `spam` to `not spam`.

**Guided Practice 1.39** What does 0.458 represent in Table 1.36? What does 0.059 represent in Table 1.37?<sup>41</sup>

<sup>41</sup>0.458 represents the proportion of spam emails that had a small number. 0.059 represents the fraction of emails with small numbers that are spam.

Ⓐ **Guided Practice 1.40** What does 0.139 at the intersection of `not spam` and `big` represent in Table 1.36? What does 0.908 represent in the Table 1.37?<sup>42</sup>

Ⓑ **Example 1.41** Data scientists use statistics to filter spam from incoming email messages. By noting specific characteristics of an email, a data scientist may be able to classify some emails as spam or not spam with high accuracy. One of those characteristics is whether the email contains no numbers, small numbers, or big numbers. Another characteristic is whether or not an email has any HTML content. A contingency table for the `spam` and `format` variables from the `email` data set are shown in Table 1.38. Recall that an HTML email is an email with the capacity for special formatting, e.g. bold text. In Table 1.38, which would be more helpful to someone hoping to classify email as spam or regular email: row or column proportions?

Such a person would be interested in how the proportion of spam changes within each email format. This corresponds to column proportions: the proportion of spam in plain text emails and the proportion of spam in HTML emails.

If we generate the column proportions, we can see that a higher fraction of plain text emails are spam ( $209/1195 = 17.5\%$ ) than compared to HTML emails ( $158/2726 = 5.8\%$ ). This information on its own is insufficient to classify an email as spam or not spam, as over 80% of plain text emails are not spam. Yet, when we carefully combine this information with many other characteristics, such as `number` and other variables, we stand a reasonable chance of being able to classify some email as spam or not spam. This is a topic we will return to in Chapter ??.

	text	HTML	Total
spam	209	158	367
not spam	986	2568	3554
Total	1195	2726	3921

Table 1.38: A contingency table for `spam` and `format`.

Example 1.41 points out that row and column proportions are not equivalent. Before settling on one form for a table, it is important to consider each to ensure that the most useful table is constructed.

Ⓐ **Guided Practice 1.42** Look back to Tables 1.36 and 1.37. Which would be more useful to someone hoping to identify spam emails using the `number` variable?<sup>43</sup>

---

<sup>42</sup>0.139 represents the fraction of non-spam email that had a big number. 0.908 represents the fraction of emails with big numbers that are non-spam emails.

<sup>43</sup>The column proportions in Table 1.37 will probably be most useful, which makes it easier to see that emails with small numbers are spam about 5.9% of the time (relatively rare). We would also see that about 27.1% of emails with no numbers are spam, and 9.2% of emails with big numbers are spam.

### 1.7.3 Segmented bar and mosaic plots

Contingency tables using row or column proportions are especially useful for examining how two categorical variables are related. Segmented bar and mosaic plots provide a way to visualize the information in these tables.

A **segmented bar plot** is a graphical display of contingency table information. For example, a segmented bar plot representing Table 1.37 is shown in Figure 1.39(a), where we have first created a bar plot using the `number` variable and then divided each group by the levels of `spam`. The column proportions of Table 1.37 have been translated into a standardized segmented bar plot in Figure 1.39(b), which is a helpful visualization of the fraction of spam emails in each level of `number`.

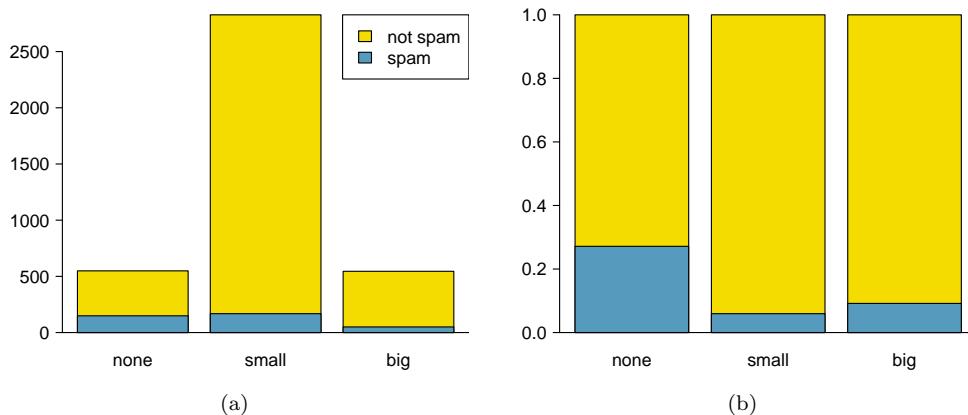


Figure 1.39: (a) Segmented bar plot for numbers found in emails, where the counts have been further broken down by `spam`. (b) Standardized version of Figure (a).

- **Example 1.43** Examine both of the segmented bar plots. Which is more useful?

Figure 1.39(a) contains more information, but Figure 1.39(b) presents the information more clearly. This second plot makes it clear that emails with no number have a relatively high rate of spam email – about 27%! On the other hand, less than 10% of email with small or big numbers are spam.

Since the proportion of spam changes across the groups in Figure 1.39(b), we can conclude the variables are dependent, which is something we were also able to discern using table proportions. Because both the `none` and `big` groups have relatively few observations compared to the `small` group, the association is more difficult to see in Figure 1.39(a).

In some other cases, a segmented bar plot that is not standardized will be more useful in communicating important information. Before settling on a particular segmented bar plot, create standardized and non-standardized forms and decide which is more effective at communicating features of the data.

A **mosaic plot** is a graphical display of contingency table information that is similar to a bar plot for one variable or a segmented bar plot when using two variables. Figure 1.40(a) shows a mosaic plot for the `number` variable. Each column represents a level of `number`, and the column widths correspond to the proportion of emails for each number type.

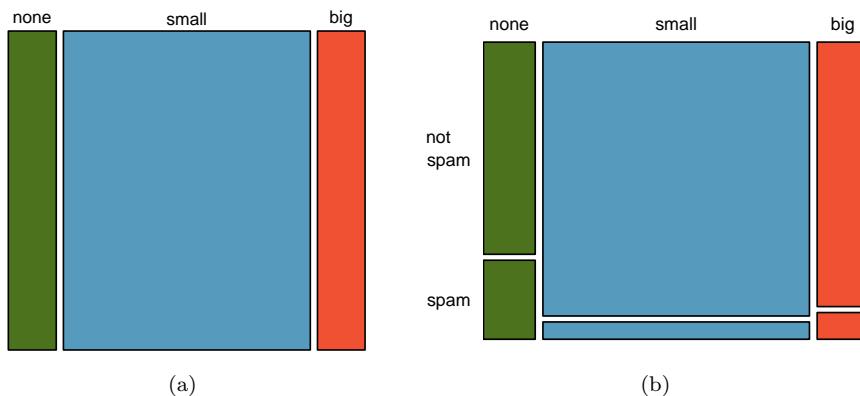


Figure 1.40: The one-variable mosaic plot for `number` and the two-variable mosaic plot for both `number` and `spam`.

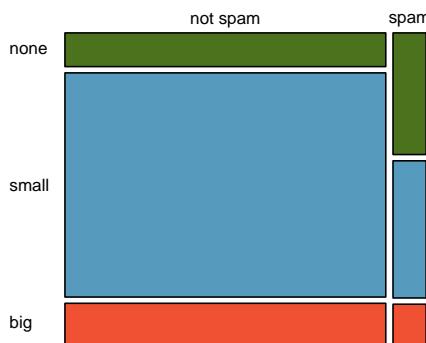


Figure 1.41: Mosaic plot where emails are grouped by the `number` variable after they've been divided into `spam` and `not spam`.

For instance, there are fewer emails with no numbers than emails with only small numbers, so the no number email column is slimmer. In general, mosaic plots use box *areas* to represent the number of observations that box represents.

This one-variable mosaic plot is further divided into pieces in Figure 1.40(b) using the `spam` variable. Each column is split proportionally according to the fraction of emails that were spam in each number category. For example, the second column, representing emails with only small numbers, was divided into emails that were spam (lower) and not spam (upper). As another example, the bottom of the third column represents spam emails that had big numbers, and the upper part of the third column represents regular emails that had big numbers. We can again use this plot to see that the `spam` and `number` variables are associated since some columns are divided in different vertical locations than others, which was the same technique used for checking an association in the standardized version of the segmented bar plot.

In a similar way, a mosaic plot representing row proportions of Table 1.33 could be constructed, as shown in Figure 1.41. However, because it is more insightful for this application to consider the fraction of spam in each category of the **number** variable, we prefer Figure 1.40(b).

### 1.7.4 The only pie chart you will see in this book

While pie charts are well known, they are not typically as useful as other charts in a data analysis. A **pie chart** is shown in Figure 1.42 alongside a bar plot. It is generally more difficult to compare group sizes in a pie chart than in a bar plot, especially when categories have nearly identical counts or proportions. In the case of the `none` and `big` categories, the difference is so slight you may be unable to distinguish any difference in group sizes for either plot!

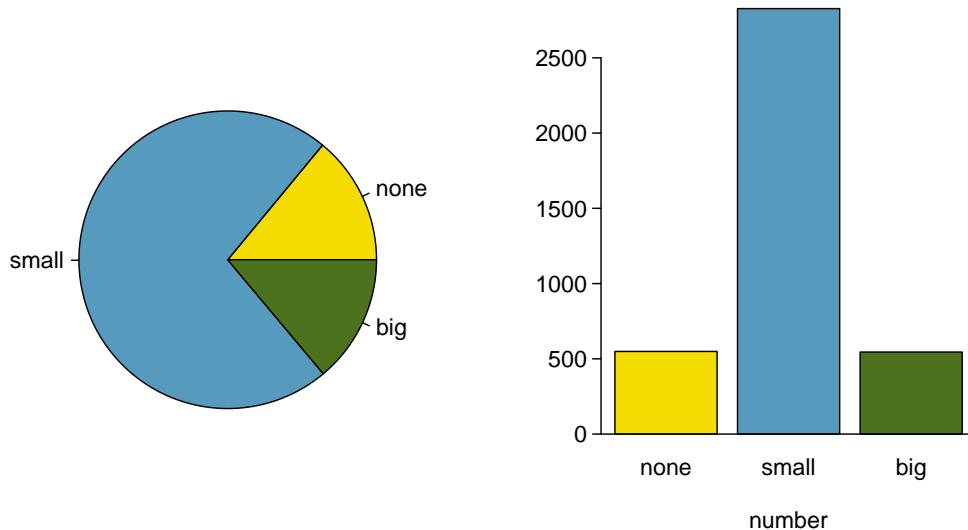


Figure 1.42: A pie chart and bar plot of `number` for the `email` data set.

### 1.7.5 Comparing numerical data across groups

Some of the more interesting investigations can be considered by examining numerical data across groups. The methods required here aren't really new. All that is required is to make a numerical plot for each group. Here two convenient methods are introduced: side-by-side box plots and hollow histograms.

We will take a look again at the `county` data set and compare the median household income for counties that gained population from 2000 to 2010 versus counties that had no gain. While we might like to make a causal connection here, remember that these are observational data and so such an interpretation would be unjustified.

There were 2,041 counties where the population increased from 2000 to 2010, and there were 1,099 counties with no gain (all but one were a loss). A random sample of 100 counties from the first group and 50 from the second group are shown in Table 1.43 to give a better sense of some of the raw data.

The **side-by-side box plot** is a traditional tool for comparing across groups. An example is shown in the left panel of Figure 1.44, where there are two box plots, one for each group, placed into one plotting window and drawn on the same scale.

Another useful plotting method uses **hollow histograms** to compare numerical data across groups. These are just the outlines of histograms of each group put on the same plot, as shown in the right panel of Figure 1.44.

population gain						no gain		
41.2	33.1	30.4	37.3	79.1	34.5	40.3	33.5	34.8
22.9	39.9	31.4	45.1	50.6	59.4	29.5	31.8	41.3
47.9	36.4	42.2	43.2	31.8	36.9	28	39.1	42.8
50.1	27.3	37.5	53.5	26.1	57.2	38.1	39.5	22.3
57.4	42.6	40.6	48.8	28.1	29.4	43.3	37.5	47.1
43.8	26	33.8	35.7	38.5	42.3	43.7	36.7	36
41.3	40.5	68.3	31	46.7	30.5	35.8	38.7	39.8
68.3	48.3	38.7	62	37.6	32.2	46	42.3	48.2
42.6	53.6	50.7	35.1	30.6	56.8	38.6	31.9	31.1
66.4	41.4	34.3	38.9	37.3	41.7	37.6	29.3	30.1
51.9	83.3	46.3	48.4	40.8	42.6	57.5	32.6	31.1
44.5	34	48.7	45.2	34.7	32.2	46.2	26.5	40.1
39.4	38.6	40	57.3	45.2	33.1	38.4	46.7	25.9
43.8	71.7	45.1	32.2	63.3	54.7	36.4	41.5	45.7
71.3	36.3	36.4	41	37	66.7	39.7	37	37.7
50.2	45.8	45.7	60.2	53.1		21.4	29.3	50.1
35.8	40.4	51.5	66.4	36.1		43.6	39.8	

Table 1.43: In this table, median household income (in \$1000s) from a random sample of 100 counties that gained population over 2000-2010 are shown on the left. Median incomes from a random sample of 50 counties that had no population gain are shown on the right.

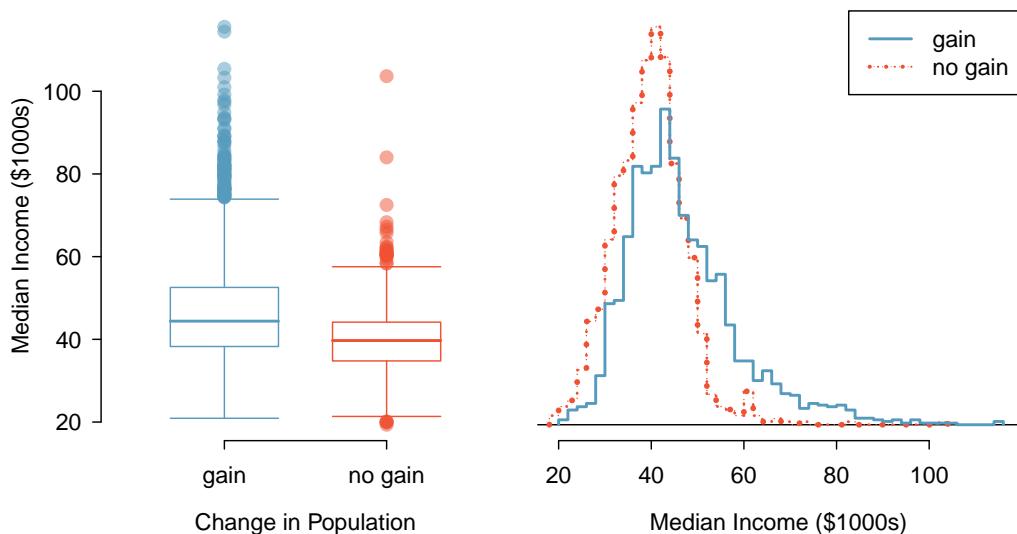


Figure 1.44: Side-by-side box plot (left panel) and hollow histograms (right panel) for `med_income`, where the counties are split by whether there was a population gain or loss from 2000 to 2010. The income data were collected between 2006 and 2010.

- Ⓐ **Guided Practice 1.44** Use the plots in Figure 1.44 to compare the incomes for counties across the two groups. What do you notice about the approximate center of each group? What do you notice about the variability between groups? Is the shape relatively consistent between groups? How many *prominent* modes are there for each group?<sup>44</sup>
- Ⓐ **Guided Practice 1.45** What components of each plot in Figure 1.44 do you find most useful?<sup>45</sup>

## 1.8 Case study: gender discrimination

(special topic)

- **Example 1.46** Suppose your professor splits the students in class into two groups: students on the left and students on the right. If  $\hat{p}_L$  and  $\hat{p}_R$  represent the proportion of students who own an Apple product on the left and right, respectively, would you be surprised if  $\hat{p}_L$  did not exactly equal  $\hat{p}_R$ ?

While the proportions would probably be close to each other, it would be unusual for them to be exactly the same. We would probably observe a small difference due to chance.

- Ⓐ **Guided Practice 1.47** If we don't think the side of the room a person sits on in class is related to whether the person owns an Apple product, what assumption are we making about the relationship between these two variables?<sup>46</sup>

### 1.8.1 Variability within data

We consider a study investigating gender discrimination in the 1970s, which is set in the context of personnel decisions within a bank.<sup>47</sup> The research question we hope to answer is, “Are females unfairly discriminated against in promotion decisions made by male managers?”

The participants in this study are 48 male bank supervisors attending a management institute at the University of North Carolina in 1972. They were asked to assume the role of the personnel director of a bank and were given a personnel file to judge whether the person should be promoted to a branch manager position. The files given to the participants were identical, except that half of them indicated the candidate was male and the other half indicated the candidate was female. These files were randomly assigned to the subjects.

---

<sup>44</sup>Answers may vary a little. The counties with population gains tend to have higher income (median of about \$45,000) versus counties without a gain (median of about \$40,000). The variability is also slightly larger for the population gain group. This is evident in the IQR, which is about 50% bigger in the *gain* group. Both distributions show slight to moderate right skew and are unimodal. There is a secondary small bump at about \$60,000 for the *no gain* group, visible in the hollow histogram plot, that seems out of place. (Looking into the data set, we would find that 8 of these 15 counties are in Alaska and Texas.) The box plots indicate there are many observations far above the median in each group, though we should anticipate that many observations will fall beyond the whiskers when using such a large data set.

<sup>45</sup>Answers will vary. The side-by-side box plots are especially useful for comparing centers and spreads, while the hollow histograms are more useful for seeing distribution shape, skew, and groups of anomalies.

<sup>46</sup>We would be assuming that these two variables are independent.

<sup>47</sup>Rosen B and Jerdee T. 1974. Influence of sex role stereotypes on personnel decisions. *Journal of Applied Psychology* 59(1):9-14.

- Ⓐ **Guided Practice 1.48** Is this an observational study or an experiment? What implications does the study type have on what can be inferred from the results?<sup>48</sup>

For each supervisor we record the gender associated with the assigned file and the promotion decision. Using the results of the study summarized in Table 1.45, we would like to evaluate if females are unfairly discriminated against in promotion decisions. In this study, a smaller proportion of females are promoted than males (0.583 versus 0.875), but it is unclear whether the difference provides *convincing evidence* that females are unfairly discriminated against.

		decision		Total
		promoted	not promoted	
gender	male	21	3	24
	female	14	10	24
	Total	35	13	48

Table 1.45: Summary results for the gender discrimination study.

- Ⓑ **Example 1.49** Statisticians are sometimes called upon to evaluate the strength of evidence. When looking at the rates of promotion for males and females in this study, what comes to mind as we try to determine whether the data show convincing evidence of a real difference?

The observed promotion rates (58.3% for females versus 87.5% for males) suggest there might be discrimination against women in promotion decisions. However, we cannot be sure if the observed difference represents discrimination or is just from random chance. Generally there is a little bit of fluctuation in sample data, and we wouldn't expect the sample proportions to be *exactly* equal, even if the truth was that the promotion decisions were independent of gender.

Example 1.49 is a reminder that the observed outcomes in the sample may not perfectly reflect the true relationships between variables in the underlying population. Table 1.45 shows there were 7 fewer promotions in the female group than in the male group, a difference in promotion rates of 29.2% ( $\frac{21}{24} - \frac{14}{24} = 0.292$ ). This difference is large, but the sample size for the study is small, making it unclear if this observed difference represents discrimination or whether it is simply due to chance. We label these two competing claims,  $H_0$  and  $H_A$ :

$H_0$ : **Independence model.** The variables `gender` and `decision` are independent. They have no relationship, and the observed difference between the proportion of males and females who were promoted, 29.2%, was due to chance.

$H_A$ : **Alternative model.** The variables `gender` and `decision` are *not* independent. The difference in promotion rates of 29.2% was not due to chance, and equally qualified females are less likely to be promoted than males.

What would it mean if the independence model, which says the variables `gender` and `decision` are unrelated, is true? It would mean each banker was going to decide whether to promote the candidate without regard to the gender indicated on the file. That is,

<sup>48</sup>The study is an experiment, as subjects were randomly assigned a male file or a female file. Since this is an experiment, the results can be used to evaluate a causal relationship between gender of a candidate and the promotion decision.

the difference in the promotion percentages was due to the way the files were randomly divided to the bankers, and the randomization just happened to give rise to a relatively large difference of 29.2%.

Consider the alternative model: bankers were influenced by which gender was listed on the personnel file. If this was true, and especially if this influence was substantial, we would expect to see some difference in the promotion rates of male and female candidates. If this gender bias was against females, we would expect a smaller fraction of promotion decisions for female personnel files relative to the male files.

We choose between these two competing claims by assessing if the data conflict so much with  $H_0$  that the independence model cannot be deemed reasonable. If this is the case, and the data support  $H_A$ , then we will reject the notion of independence and conclude there was discrimination.

### 1.8.2 Simulating the study

Table 1.45 shows that 35 bank supervisors recommended promotion and 13 did not. Now, suppose the bankers' decisions were independent of gender. Then, if we conducted the experiment again with a different random arrangement of files, differences in promotion rates would be based only on random fluctuation. We can actually perform this **randomization**, which simulates what would have happened if the bankers' decisions had been independent of gender but we had distributed the files differently.

In this **simulation**, we thoroughly shuffle 48 personnel files, 24 labeled `male_sim` and 24 labeled `female_sim`, and deal these files into two stacks. We will deal 35 files into the first stack, which will represent the 35 supervisors who recommended promotion. The second stack will have 13 files, and it will represent the 13 supervisors who recommended against promotion. Then, as we did with the original data, we tabulate the results and determine the fraction of `male_sim` and `female_sim` who were promoted. The randomization of files in this simulation is independent of the promotion decisions, which means any difference in the two fractions is entirely due to chance. Table 1.46 show the results of such a simulation.

		decision		Total
		promoted	not promoted	
gender_sim	male_sim	18	6	24
	female_sim	17	7	24
	Total	35	13	48

Table 1.46: Simulation results, where any difference in promotion rates between `male_sim` and `female_sim` is purely due to chance.

- **Guided Practice 1.50** What is the difference in promotion rates between the two simulated groups in Table 1.46? How does this compare to the observed 29.2% in the actual groups?<sup>49</sup>

---

<sup>49</sup> $18/24 - 17/24 = 0.042$  or about 4.2% in favor of the men. This difference due to chance is much smaller than the difference observed in the actual groups.

### 1.8.3 Checking for independence

We computed one possible difference under the independence model in Guided Practice 1.50, which represents one difference due to chance. While in this first simulation, we physically dealt out files, it is more efficient to perform this simulation using a computer. Repeating the simulation on a computer, we get another difference due to chance: -0.042. And another: 0.208. And so on until we repeat the simulation enough times that we have a good idea of what represents the *distribution of differences from chance alone*. Figure 1.47 shows a plot of the differences found from 100 simulations, where each dot represents a simulated difference between the proportions of male and female files that were recommended for promotion.

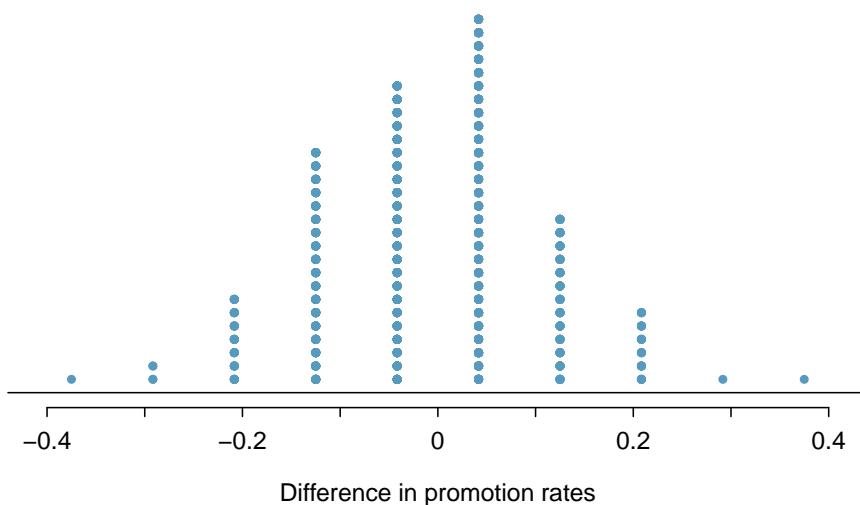


Figure 1.47: A stacked dot plot of differences from 100 simulations produced under the independence model,  $H_0$ , where `gender_sim` and `decision` are independent. Two of the 100 simulations had a difference of at least 29.2%, the difference observed in the study.

Note that the distribution of these simulated differences is centered around 0. We simulated these differences assuming that the independence model was true, and under this condition, we expect the difference to be zero with some random fluctuation. We would generally be surprised to see a difference of *exactly* 0: sometimes, just by chance, the difference is higher than 0, and other times it is lower than zero.

- **Example 1.51** How often would you observe a difference of at least 29.2% (0.292) according to Figure 1.47? Often, sometimes, rarely, or never?

---

It appears that a difference of at least 29.2% due to chance alone would only happen about 2% of the time according to Figure 1.47. Such a low probability indicates a rare event.

The difference of 29.2% being a rare event suggests two possible interpretations of the results of the study:

$H_0$  **Independence model.** Gender has no effect on promotion decision, and we observed a difference that would only happen rarely.

$H_A$  **Alternative model.** Gender has an effect on promotion decision, and what we observed was actually due to equally qualified women being discriminated against in promotion decisions, which explains the large difference of 29.2%.

Based on the simulations, we have two options. (1) We conclude that the study results do not provide strong evidence against the independence model. That is, we do not have sufficiently strong evidence to conclude there was gender discrimination. (2) We conclude the evidence is sufficiently strong to reject  $H_0$  and assert that there was gender discrimination. When we conduct formal studies, usually we reject the notion that we just happened to observe a rare event.<sup>50</sup> So in this case, we reject the independence model in favor of the alternative. That is, we are concluding the data provide strong evidence of gender discrimination against women by the supervisors.

One field of statistics, statistical inference, is built on evaluating whether such differences are due to chance. In statistical inference, statisticians evaluate which model is most reasonable given the data. Errors do occur, just like rare events, and we might choose the wrong model. While we do not always choose correctly, statistical inference gives us tools to control and evaluate how often these errors occur. In Chapter 4, we give a formal introduction to the problem of model selection. We spend the next two chapters building a foundation of probability and theory necessary to make that discussion rigorous.

---

<sup>50</sup>This reasoning does not generally extend to anecdotal observations. Each of us observes incredibly rare events every day, events we could not possibly hope to predict. However, in the non-rigorous setting of anecdotal evidence, almost anything may appear to be a rare event, so the idea of looking for rare events in day-to-day activities is treacherous. For example, we might look at the lottery: there was only a 1 in 176 million chance that the Mega Millions numbers for the largest jackpot in history (March 30, 2012) would be (2, 4, 23, 38, 46) with a Mega ball of (23), but nonetheless those numbers came up! However, no matter what numbers had turned up, they would have had the same incredibly rare odds. That is, *any set of numbers we could have observed would ultimately be incredibly rare*. This type of situation is typical of our daily lives: each possible event in itself seems incredibly rare, but if we consider every alternative, those outcomes are also incredibly rare. We should be cautious not to misinterpret such anecdotal evidence.

## 1.9 Exercises

### 1.9.1 Case study: using stents to prevent strokes

**1.1 Migraine and acupuncture.** A migraine is a particularly painful type of headache, which patients sometimes wish to treat with acupuncture. To determine whether acupuncture relieves migraine pain, researchers conducted a randomized controlled study where 89 females diagnosed with migraine headaches were randomly assigned to one of two groups: treatment or control. 43 patients in the treatment group received acupuncture that is specifically designed to treat migraines. 46 patients in the control group received placebo acupuncture (needle insertion at non-acupoint locations). 24 hours after patients received acupuncture, they were asked if they were pain free. Results are summarized in the contingency table below.<sup>51</sup>

		Pain free		Total
		Yes	No	
Group	Treatment	10	33	43
	Control	2	44	46
	Total	12	77	89



Figure from the original paper displaying the appropriate area (M) versus the inappropriate area (S) used in the treatment of migraine attacks.

- (a) What percent of patients in the treatment group were pain free 24 hours after receiving acupuncture? What percent in the control group?
- (b) At first glance, does acupuncture appear to be an effective treatment for migraines? Explain your reasoning.
- (c) Do the data provide convincing evidence that there is a real pain reduction for those patients in the treatment group? Or do you think that the observed difference might just be due to chance?

**1.2 Sinusitis and antibiotics.** Researchers studying the effect of antibiotic treatment for acute sinusitis compared to symptomatic treatments randomly assigned 166 adults diagnosed with acute sinusitis to one of two groups: treatment or control. Study participants received either a 10-day course of amoxicillin (an antibiotic) or a placebo similar in appearance and taste. The placebo consisted of symptomatic treatments such as acetaminophen, nasal decongestants, etc. At the end of the 10-day period patients were asked if they experienced significant improvement in symptoms. The distribution of responses is summarized below.<sup>52</sup>

		Self-reported significant improvement in symptoms		Total
		Yes	No	
Group	Treatment	66	19	85
	Control	65	16	81
	Total	131	35	166

- (a) What percent of patients in the treatment group experienced a significant improvement in symptoms? What percent in the control group?
- (b) Based on your findings in part (a), which treatment appears to be more effective for sinusitis?
- (c) Do the data provide convincing evidence that there is a difference in the improvement rates of sinusitis symptoms? Or do you think that the observed difference might just be due to chance?

<sup>51</sup>G. Allais et al. “Ear acupuncture in the treatment of migraine attacks: a randomized trial on the efficacy of appropriate versus inappropriate acupoints”. In: *Neurological Sci.* 32.1 (2011), pp. 173–175.

<sup>52</sup>J.M. Garbutt et al. “Amoxicillin for Acute Rhinosinusitis: A Randomized Controlled Trial”. In: *JAMA: The Journal of the American Medical Association* 307.7 (2012), pp. 685–692.

## 1.9.2 Data basics

**1.3 Air pollution and birth outcomes, study components.** Researchers collected data to examine the relationship between air pollutants and preterm births in Southern California. During the study air pollution levels were measured by air quality monitoring stations. Specifically, levels of carbon monoxide were recorded in parts per million, nitrogen dioxide and ozone in parts per hundred million, and coarse particulate matter ( $\text{PM}_{10}$ ) in  $\mu\text{g}/\text{m}^3$ . Length of gestation data were collected on 143,196 births between the years 1989 and 1993, and air pollution exposure during gestation was calculated for each birth. The analysis suggested that increased ambient  $\text{PM}_{10}$  and, to a lesser degree, CO concentrations may be associated with the occurrence of preterm births.<sup>53</sup> In this study, identify

- (a) the cases,
- (b) the variables and their types, and
- (c) the main research question.

**1.4 Buteyko method, study components.** The Buteyko method is a shallow breathing technique developed by Konstantin Buteyko, a Russian doctor, in 1952. Anecdotal evidence suggests that the Buteyko method can reduce asthma symptoms and improve quality of life. In a scientific study to determine the effectiveness of this method, researchers recruited 600 asthma patients aged 18-69 who relied on medication for asthma treatment. These patients were split into two research groups: one practiced the Buteyko method and the other did not. Patients were scored on quality of life, activity, asthma symptoms, and medication reduction on a scale from 0 to 10. On average, the participants in the Buteyko group experienced a significant reduction in asthma symptoms and an improvement in quality of life.<sup>54</sup> In this study, identify

- (a) the cases,
- (b) the variables and their types, and
- (c) the main research question.

**1.5 Cheaters, study components.** Researchers studying the relationship between honesty, age and self-control conducted an experiment on 160 children between the ages of 5 and 15. Participants reported their age, sex, and whether they were an only child or not. The researchers asked each child to toss a fair coin in private and to record the outcome (white or black) on a paper sheet, and said they would only reward children who report white. Half the students were explicitly told not to cheat and the others were not given any explicit instructions. In the no instruction group probability of cheating was found to be uniform across groups based on child's characteristics. In the group that was explicitly told to not cheat, girls were less likely to cheat, and while rate of cheating didn't vary by age for boys, it decreased with age for girls.<sup>55</sup> In this study, identify

- (a) the cases,
- (b) the variables and their types, and
- (c) the main research question.

---

<sup>53</sup>B. Ritz et al. "Effect of air pollution on preterm birth among children born in Southern California between 1989 and 1993". In: *Epidemiology* 11.5 (2000), pp. 502–511.

<sup>54</sup>J. McGowan. "Health Education: Does the Buteyko Institute Method make a difference?" In: *Thorax* 58 (2003).

<sup>55</sup>Alessandro Bucciol and Marco Piovesan. "Luck or cheating? A field experiment on honesty with children". In: *Journal of Economic Psychology* 32.1 (2011), pp. 73–78.

**1.6 Stealers, study components.** In a study of the relationship between socio-economic class and unethical behavior, 129 University of California undergraduates at Berkeley were asked to identify themselves as having low or high social-class by comparing themselves to others with the most (least) money, most (least) education, and most (least) respected jobs. They were also presented with a jar of individually wrapped candies and informed that the candies were for children in a nearby laboratory, but that they could take some if they wanted. After completing some unrelated tasks, participants reported the number of candies they had taken. It was found that those who were identified as upper-class took more candy than others.<sup>56</sup> In this study, identify

- (a) the cases,
- (b) the variables and their types, and
- (c) the main research question.

**1.7 Fisher's irises.** Sir Ronald Aylmer Fisher was an English statistician, evolutionary biologist, and geneticist who worked on a data set that contained sepal length and width, and petal length and width from three species of iris flowers (*setosa*, *versicolor* and *virginica*). There were 50 flowers from each species in the data set.<sup>57</sup>

- (a) How many cases were included in the data?
- (b) How many numerical variables are included in the data? Indicate what they are, and if they are continuous or discrete.
- (c) How many categorical variables are included in the data, and what are they? List the corresponding levels (categories).



Photo by Ryan Claussen  
(<http://flic.kr/p/6QTcuX>)  
CC BY-SA 2.0 license

**1.8 Smoking habits of UK residents.** A survey was conducted to study the smoking habits of UK residents. Below is a data matrix displaying a portion of the data collected in this survey. Note that “£” stands for British Pounds Sterling, “cig” stands for cigarettes, and “N/A” refers to a missing component of the data.<sup>58</sup>

	sex	age	marital	grossIncome	smoke	amtWeekends	amtWeekdays
1	Female	42	Single	Under £2,600	Yes	12 cig/day	12 cig/day
2	Male	44	Single	£10,400 to £15,600	No	N/A	N/A
3	Male	53	Married	Above £36,400	Yes	6 cig/day	6 cig/day
:	:	:	:	:	:	:	:
1691	Male	40	Single	£2,600 to £5,200	Yes	8 cig/day	8 cig/day

- (a) What does each row of the data matrix represent?
- (b) How many participants were included in the survey?
- (c) Indicate whether each variable in the study is numerical or categorical. If numerical, identify as continuous or discrete. If categorical, indicate if the variable is ordinal.

<sup>56</sup>P.K. Piff et al. “Higher social class predicts increased unethical behavior”. In: *Proceedings of the National Academy of Sciences* (2012).

<sup>57</sup>R.A Fisher. “The Use of Multiple Measurements in Taxonomic Problems”. In: *Annals of Eugenics* 7 (1936), pp. 179–188.

<sup>58</sup>National STEM Centre, Large Datasets from stats4schools.

### 1.9.3 Overview of data collection principles

**1.9 Air pollution and birth outcomes, scope of inference.** Exercise 1.3 introduces a study where researchers collected data to examine the relationship between air pollutants and preterm births in Southern California. During the study air pollution levels were measured by air quality monitoring stations. Length of gestation data were collected on 143,196 births between the years 1989 and 1993, and air pollution exposure during gestation was calculated for each birth.

- (a) Identify the population of interest and the sample in this study.
- (b) Comment on whether or not the results of the study can be generalized to the population, and if the findings of the study can be used to establish causal relationships.

**1.10 Cheaters, scope of inference.** Exercise 1.5 introduces a study where researchers studying the relationship between honesty, age, and self-control conducted an experiment on 160 children between the ages of 5 and 15. The researchers asked each child to toss a fair coin in private and to record the outcome (white or black) on a paper sheet, and said they would only reward children who report white. Half the students were explicitly told not to cheat and the others were not given any explicit instructions. Differences were observed in the cheating rates in the instruction and no instruction groups, as well as some differences across children's characteristics within each group.

- (a) Identify the population of interest and the sample in this study.
- (b) Comment on whether or not the results of the study can be generalized to the population, and if the findings of the study can be used to establish causal relationships.

**1.11 Buteyko method, scope of inference.** Exercise 1.4 introduces a study on using the Buteyko shallow breathing technique to reduce asthma symptoms and improve quality of life. As part of this study 600 asthma patients aged 18-69 who relied on medication for asthma treatment were recruited and randomly assigned to two groups: one practiced the Buteyko method and the other did not. Those in the Buteyko group experienced, on average, a significant reduction in asthma symptoms and an improvement in quality of life.

- (a) Identify the population of interest and the sample in this study.
- (b) Comment on whether or not the results of the study can be generalized to the population, and if the findings of the study can be used to establish causal relationships.

**1.12 Stealers, scope of inference.** Exercise 1.6 introduces a study on the relationship between socio-economic class and unethical behavior. As part of this study 129 University of California Berkeley undergraduates were asked to identify themselves as having low or high social-class by comparing themselves to others with the most (least) money, most (least) education, and most (least) respected jobs. They were also presented with a jar of individually wrapped candies and informed that the candies were for children in a nearby laboratory, but that they could take some if they wanted. After completing some unrelated tasks, participants reported the number of candies they had taken. It was found that those who were identified as upper-class took more candy than others.

- (a) Identify the population of interest and the sample in this study.
- (b) Comment on whether or not the results of the study can be generalized to the population, and if the findings of the study can be used to establish causal relationships.

**1.13 Relaxing after work.** The 2010 General Social Survey asked the question, “After an average work day, about how many hours do you have to relax or pursue activities that you enjoy?” to a random sample of 1,155 Americans. The average relaxing time was found to be 1.65 hours. Determine which of the following is an observation, a variable, a sample statistic, or a population parameter.

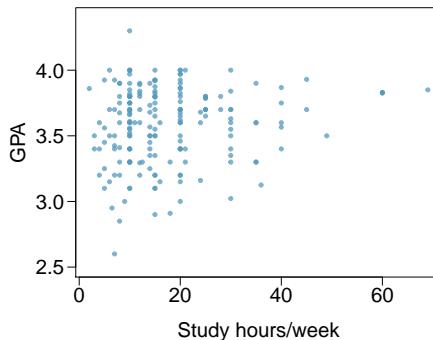
- (a) An American in the sample.
- (b) Number of hours spent relaxing after an average work day.
- (c) 1.65.
- (d) Average number of hours all Americans spend relaxing after an average work day.

**1.14 Cats on YouTube.** Suppose you want to estimate the percentage of videos on YouTube that are cat videos. It is impossible for you to watch all videos on YouTube so you use a random video picker to select 1000 videos for you. You find that 2% of these videos are cat videos. Determine which of the following is an observation, a variable, a sample statistic, or a population parameter.

- (a) Percentage of all videos on YouTube that are cat videos.
- (b) 2%.
- (c) A video in your sample.
- (d) Whether or not a video is a cat video.

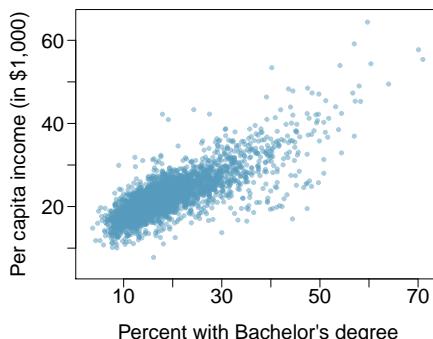
**1.15 GPA and study hours.** A survey was conducted on 193 Duke University undergraduates who took an introductory statistics course in 2012. Among many other questions, this survey asked them about their GPA, which can range between 0 and 4 points, and the number of hours they spent studying per week. The scatterplot below displays the relationship between these two variables.

- (a) What is the explanatory variable and what is the response variable?
- (b) Describe the relationship between the two variables. Make sure to discuss unusual observations, if any.
- (c) Is this an experiment or an observational study?
- (d) Can we conclude that studying longer hours leads to higher GPAs?



**1.16 Income and education in US counties.** The scatterplot below shows the relationship between per capita income (in thousands of dollars) and percent of population with a bachelor's degree in 3,143 counties in the US in 2010.

- (a) What are the explanatory and response variables?
- (b) Describe the relationship between the two variables. Make sure to discuss unusual observations, if any.
- (c) Can we conclude that having a bachelor's degree increases one's income?



### 1.9.4 Observational studies and sampling strategies

**1.17 Course satisfaction across sections.** A large college class has 160 students. All 160 students attend the lectures together, but the students are divided into 4 groups, each of 40 students, for lab sections administered by different teaching assistants. The professor wants to conduct a survey about how satisfied the students are with the course, and he believes that the lab section a student is in might affect the student's overall satisfaction with the course.

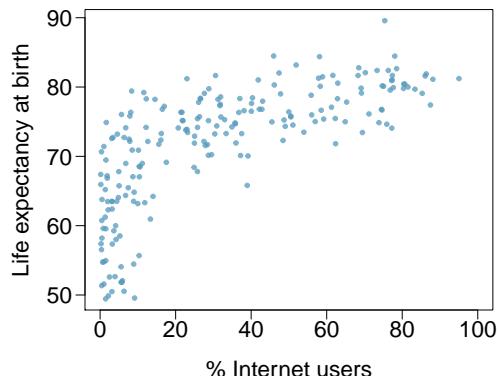
- (a) What type of study is this?
- (b) Suggest a sampling strategy for carrying out this study.

**1.18 Housing proposal across dorms.** On a large college campus first-year students and sophomores live in dorms located on the eastern part of the campus and juniors and seniors live in dorms located on the western part of the campus. Suppose you want to collect student opinions on a new housing structure the college administration is proposing and you want to make sure your survey equally represents opinions from students from all years.

- (a) What type of study is this?
- (b) Suggest a sampling strategy for carrying out this study.

**1.19 Internet use and life expectancy.** The following scatterplot was created as part of a study evaluating the relationship between estimated life expectancy at birth (as of 2014) and percentage of internet users (as of 2009) in 208 countries for which such data were available.<sup>59</sup>

- (a) Describe the relationship between life expectancy and percentage of internet users.
- (b) What type of study is this?
- (c) State a possible confounding variable that might explain this relationship and describe its potential effect.



**1.20 Stressed out, Part I.** A study that surveyed a random sample of otherwise healthy high school students found that they are more likely to get muscle cramps when they are stressed. The study also noted that students drink more coffee and sleep less when they are stressed.

- (a) What type of study is this?
- (b) Can this study be used to conclude a causal relationship between increased stress and muscle cramps?
- (c) State possible confounding variables that might explain the observed relationship between increased stress and muscle cramps.

**1.21 Evaluate sampling methods.** A university wants to determine what fraction of its undergraduate student body support a new \$25 annual fee to improve the student union. For each proposed method below, indicate whether the method is reasonable or not.

- (a) Survey a simple random sample of 500 students.
- (b) Stratify students by their field of study, then sample 10% of students from each stratum.
- (c) Cluster students by their ages (e.g. 18 years old in one cluster, 19 years old in one cluster, etc.), then randomly sample three clusters and survey all students in those clusters.

---

<sup>59</sup>CIA Factbook, Country Comparisons, 2014.

**1.22 Random digit dialing.** The Gallup Poll uses a procedure called random digit dialing, which creates phone numbers based on a list of all area codes in America in conjunction with the associated number of residential households in each area code. Give a possible reason the Gallup Poll chooses to use random digit dialing instead of picking phone numbers from the phone book.

**1.23 Haters are gonna hate, study confirms.** A study published in the *Journal of Personality and Social Psychology* asked a group of 200 randomly sampled men and women to evaluate how they felt about various subjects, such as camping, health care, architecture, taxidermy, crossword puzzles, and Japan in order to measure their dispositional attitude towards mostly independent stimuli. Then, they presented the participants with information about a new product: a microwave oven. This microwave oven does not exist, but the participants didn't know this, and were given three positive and three negative fake reviews. People who reacted positively to the subjects on the dispositional attitude measurement also tended to react positively to the microwave oven, and those who reacted negatively also tended to react negatively to it. Researchers concluded that "some people tend to like things, whereas others tend to dislike things, and a more thorough understanding of this tendency will lead to a more thorough understanding of the psychology of attitudes."<sup>60</sup>

- (a) What are the cases?
- (b) What is (are) the response variable(s) in this study?
- (c) What is (are) the explanatory variable(s) in this study?
- (d) Does the study employ random sampling?
- (e) Is this an observational study or an experiment? Explain your reasoning.
- (f) Can we establish a causal link between the explanatory and response variables?
- (g) Can the results of the study be generalized to the population at large?

**1.24 Family size.** Suppose we want to estimate household size, where a "household" is defined as people living together in the same dwelling, and sharing living accommodations. If we select students at random at an elementary school and ask them what their family size is, will this be a good measure of household size? Or will our average be biased? If so, will it overestimate or underestimate the true value?

**1.25 Flawed reasoning.** Identify the flaw(s) in reasoning in the following scenarios. Explain what the individuals in the study should have done differently if they wanted to make such strong conclusions.

- (a) Students at an elementary school are given a questionnaire that they are asked to return after their parents have completed it. One of the questions asked is, "Do you find that your work schedule makes it difficult for you to spend time with your kids after school?" Of the parents who replied, 85% said "no". Based on these results, the school officials conclude that a great majority of the parents have no difficulty spending time with their kids after school.
- (b) A survey is conducted on a simple random sample of 1,000 women who recently gave birth, asking them about whether or not they smoked during pregnancy. A follow-up survey asking if the children have respiratory problems is conducted 3 years later, however, only 567 of these women are reached at the same address. The researcher reports that these 567 women are representative of all mothers.
- (c) An orthopedist administers a questionnaire to 30 of his patients who do not have any joint problems and finds that 20 of them regularly go running. He concludes that running decreases the risk of joint problems.

---

<sup>60</sup>Justin Hepler and Dolores Albarracín. "Attitudes without objects - Evidence for a dispositional attitude, its measurement, and its consequences". In: *Journal of personality and social psychology* 104.6 (2013), p. 1060.

**1.26 City council survey.** A city council has requested a household survey be conducted in a suburban area of their city. The area is broken into many distinct and unique neighborhoods, some including large homes, some with only apartments, and others a diverse mixture of housing structures. Identify the sampling methods described below, and comment on whether or not you think they would be effective in this setting.

- (a) Randomly sample 50 households from the city.
- (b) Divide the city into neighborhoods, and sample 20 households from each neighborhood.
- (c) Divide the city into neighborhoods, randomly sample 10 neighborhoods, and sample all households from those neighborhoods.
- (d) Divide the city into neighborhoods, randomly sample 10 neighborhoods, and then randomly sample 20 households from those neighborhoods.
- (e) Sample the 200 households closest to the city council offices.

**1.27 Sampling strategies.** A statistics student who is curious about the relationship between the amount of time students spend on social networking sites and their performance at school decides to conduct a survey. Various research strategies for collecting data are described below. In each, name the sampling method proposed and any bias you might expect.

- (a) He randomly samples 40 students from the study's population, gives them the survey, asks them to fill it out and bring it back the next day.
- (b) He gives out the survey only to his friends, making sure each one of them fills out the survey.
- (c) He posts a link to an online survey on Facebook and asks his friends to fill out the survey.
- (d) He randomly samples 5 classes and asks a random sample of students from those classes to fill out the survey.

**1.28 Reading the paper.** Below are excerpts from two articles published in the *NY Times*:

- (a) An article titled *Risks: Smokers Found More Prone to Dementia* states the following:<sup>61</sup>  
“Researchers analyzed data from 23,123 health plan members who participated in a voluntary exam and health behavior survey from 1978 to 1985, when they were 50-60 years old. 23 years later, about 25% of the group had dementia, including 1,136 with Alzheimer’s disease and 416 with vascular dementia. After adjusting for other factors, the researchers concluded that pack-a-day smokers were 37% more likely than nonsmokers to develop dementia, and the risks went up with increased smoking; 44% for one to two packs a day; and twice the risk for more than two packs.”

Based on this study, can we conclude that smoking causes dementia later in life? Explain your reasoning.

- (b) Another article titled *The School Bully Is Sleepy* states the following:<sup>62</sup>  
“The University of Michigan study, collected survey data from parents on each child’s sleep habits and asked both parents and teachers to assess behavioral concerns. About a third of the students studied were identified by parents or teachers as having problems with disruptive behavior or bullying. The researchers found that children who had behavioral issues and those who were identified as bullies were twice as likely to have shown symptoms of sleep disorders.”

A friend of yours who read the article says, “The study shows that sleep disorders lead to bullying in school children.” Is this statement justified? If not, how best can you describe the conclusion that can be drawn from this study?

---

<sup>61</sup>R.C. Rabin. “Risks: Smokers Found More Prone to Dementia”. In: *New York Times* (2010).

<sup>62</sup>T. Parker-Pope. “The School Bully Is Sleepy”. In: *New York Times* (2011).

**1.29 Shyness on Facebook.** Given the anonymity afforded to individuals in online interactions, researchers hypothesized that shy individuals might have more favorable attitudes toward Facebook, and that shyness might be positively correlated with time spent on Facebook. They also hypothesized that shy individuals might have fewer Facebook “friends” as they tend to have fewer friends than non-shy individuals have in the offline world. 103 undergraduate students at an Ontario university were surveyed via online questionnaires. The study states “Participants were recruited through the university’s psychology participation pool. After indicating an interest in the study, participants were sent an e-mail containing the study’s URL.” Are the results of this study generalizable to the population of all Facebook users?<sup>63</sup>

### 1.9.5 Experiments

**1.30 Stressed out, Part II.** In a study evaluating the relationship between stress and muscle cramps, half the subjects are randomly assigned to be exposed to increased stress by being placed into an elevator that falls rapidly and stops abruptly and the other half are left at no or baseline stress.

- (a) What type of study is this?
- (b) Can this study be used to conclude a causal relationship between increased stress and muscle cramps?

**1.31 Light and exam performance.** A study is designed to test the effect of light level on exam performance of students. The researcher believes that light levels might have different effects on males and females, so wants to make sure both are equally represented in each treatment. The treatments are fluorescent overhead lighting, yellow overhead lighting, no overhead lighting (only desk lamps).

- (a) What is the response variable?
- (b) What is the explanatory variable? What are its levels?
- (c) What is the blocking variable? What are its levels?

**1.32 Vitamin supplements.** In order to assess the effectiveness of taking large doses of vitamin C in reducing the duration of the common cold, researchers recruited 400 healthy volunteers from staff and students at a university. A quarter of the patients were assigned a placebo, and the rest were evenly divided between 1g Vitamin C, 3g Vitamin C, or 3g Vitamin C plus additives to be taken at onset of a cold for the following two days. All tablets had identical appearance and packaging. The nurses who handed the prescribed pills to the patients knew which patient received which treatment, but the researchers assessing the patients when they were sick did not. No significant differences were observed in any measure of cold duration or severity between the four medication groups, and the placebo group had the shortest duration of symptoms.<sup>64</sup>

- (a) Was this an experiment or an observational study? Why?
- (b) What are the explanatory and response variables in this study?
- (c) Were the patients blinded to their treatment?
- (d) Was this study double-blind?
- (e) Participants are ultimately able to choose whether or not to use the pills prescribed to them. We might expect that not all of them will adhere and take their pills. Does this introduce a confounding variable to the study? Explain your reasoning.

---

<sup>63</sup>E.S. Orr et al. “The influence of shyness on the use of Facebook in an undergraduate sample”. In: *CyberPsychology & Behavior* 12.3 (2009), pp. 337–340.

<sup>64</sup>C. Audera et al. “Mega-dose vitamin C in treatment of the common cold: a randomised controlled trial”. In: *Medical Journal of Australia* 175.7 (2001), pp. 359–362.

**1.33 Light, noise, and exam performance.** A study is designed to test the effect of light level and noise level on exam performance of students. The researcher believes that light and noise levels might have different effects on males and females, so wants to make sure both are equally represented in each treatment. The light treatments considered are fluorescent overhead lighting, yellow overhead lighting, no overhead lighting (only desk lamps). The noise treatments considered are no noise, construction noise, and human chatter noise.

- (a) What is the response variable?
- (b) How many factors are considered in this study? Identify them, and describe their levels.
- (c) What is the role of the sex variable in this study?

**1.34 Music and learning.** You would like to conduct an experiment in class to see if students learn better if they study without any music, with music that has no lyrics (instrumental), or with music that has lyrics. Briefly outline a design for this study.

**1.35 Soda preference.** You would like to conduct an experiment in class to see if your classmates prefer the taste of regular Coke or Diet Coke. Briefly outline a design for this study.

**1.36 Exercise and mental health.** A researcher is interested in the effects of exercise on mental health and he proposes the following study: Use stratified random sampling to ensure representative proportions of 18-30, 31-40 and 41- 55 year olds from the population. Next, randomly assign half the subjects from each age group to exercise twice a week, and instruct the rest not to exercise. Conduct a mental health exam at the beginning and at the end of the study, and compare the results.

- (a) What type of study is this?
- (b) What are the treatment and control groups in this study?
- (c) Does this study make use of blocking? If so, what is the blocking variable?
- (d) Does this study make use of blinding?
- (e) Comment on whether or not the results of the study can be used to establish a causal relationship between exercise and mental health, and indicate whether or not the conclusions can be generalized to the population at large.
- (f) Suppose you are given the task of determining if this proposed study should get funding. Would you have any reservations about the study proposal?

**1.37 Chia seeds and weight loss.** Chia Pets – those terra-cotta figurines that sprout fuzzy green hair – made the chia plant a household name. But chia has gained an entirely new reputation as a diet supplement. In one 2009 study, a team of researchers recruited 38 men and divided them randomly into two groups: treatment or control. They also recruited 38 women, and they randomly placed half of these participants into the treatment group and the other half into the control group. One group was given 25 grams of chia seeds twice a day, and the other was given a placebo. The subjects volunteered to be a part of the study. After 12 weeks, the scientists found no significant difference between the groups in appetite or weight loss.<sup>65</sup>

- (a) What type of study is this?
- (b) What are the experimental and control treatments in this study?
- (c) Has blocking been used in this study? If so, what is the blocking variable?
- (d) Has blinding been used in this study?
- (e) Comment on whether or not we can make a causal statement, and indicate whether or not we can generalize the conclusion to the population at large.

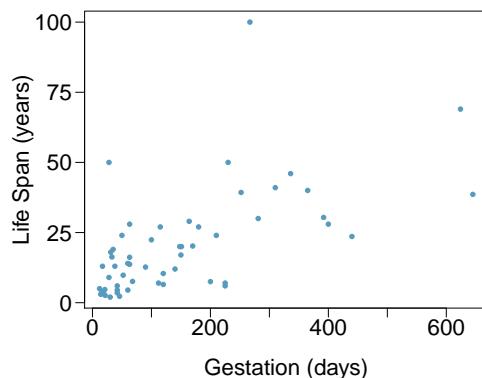
---

<sup>65</sup>D.C. Nieman et al. “Chia seed does not promote weight loss or alter disease risk factors in overweight adults”. In: *Nutrition Research* 29.6 (2009), pp. 414–418.

### 1.9.6 Examining numerical data

**1.38 Mammal life spans.** Data were collected on life spans (in years) and gestation lengths (in days) for 62 mammals. A scatterplot of life span versus length of gestation is shown below.<sup>66</sup>

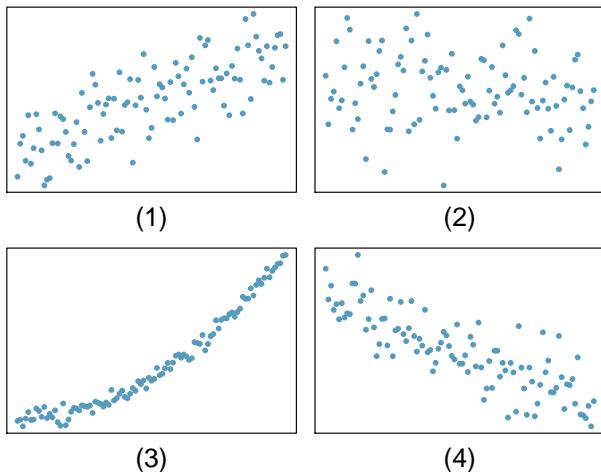
- (a) What type of an association is apparent between life span and length of gestation?
- (b) What type of an association would you expect to see if the axes of the plot were reversed, i.e. if we plotted length of gestation versus life span?
- (c) Are life span and length of gestation independent? Explain your reasoning.



**1.39 Associations.** Indicate which of the plots show a

- (a) positive association
- (b) negative association
- (c) no association

Also determine if the positive and negative associations are linear or nonlinear. Each part may refer to more than one plot.



**1.40 Office productivity.** Office productivity is relatively low when the employees feel no stress about their work or job security. However, high levels of stress can also lead to reduced employee productivity. Sketch a plot to represent the relationship between stress and productivity.

**1.41 Reproducing bacteria.** Suppose that there is only sufficient space and nutrients to support one million bacterial cells in a petri dish. You place a few bacterial cells in this petri dish, allow them to reproduce freely, and record the number of bacterial cells in the dish over time. Sketch a plot representing the relationship between number of bacterial cells and time.

**1.42 Sleeping in college.** A recent article in a college newspaper stated that college students get an average of 5.5 hrs of sleep each night. A student who was skeptical about this value decided to conduct a survey by randomly sampling 25 students. On average, the sampled students slept 6.25 hours per night. Identify which value represents the sample mean and which value represents the claimed population mean.

---

<sup>66</sup>T. Allison and D.V. Cicchetti. "Sleep in mammals: ecological and constitutional correlates". In: *Arch. Hydrobiol.* 75 (1975), p. 442.

**1.43 Parameters and statistics.** Identify which value represents the sample mean and which value represents the claimed population mean.

- (a) American households spent an average of about \$52 in 2007 on Halloween merchandise such as costumes, decorations and candy. To see if this number had changed, researchers conducted a new survey in 2008 before industry numbers were reported. The survey included 1,500 households and found that average Halloween spending was \$58 per household.
- (b) The average GPA of students in 2001 at a private university was 3.37. A survey on a sample of 203 students from this university yielded an average GPA of 3.59 in Spring semester of 2012.

**1.44 Make-up exam.** In a class of 25 students, 24 of them took an exam in class and 1 student took a make-up exam the following day. The professor graded the first batch of 24 exams and found an average score of 74 points with a standard deviation of 8.9 points. The student who took the make-up the following day scored 64 points on the exam.

- (a) Does the new student's score increase or decrease the average score?
- (b) What is the new average?
- (c) Does the new student's score increase or decrease the standard deviation of the scores?

**1.45 Days off at a mining plant.** Workers at a particular mining site receive an average of 35 days paid vacation, which is lower than the national average. The manager of this plant is under pressure from a local union to increase the amount of paid time off. However, he does not want to give more days off to the workers because that would be costly. Instead he decides he should fire 10 employees in such a way as to raise the average number of days off that are reported by his employees. In order to achieve this goal, should he fire employees who have the most number of days off, least number of days off, or those who have about the average number of days off?

**1.46 Medians and IQRs.** For each part, compare distributions (1) and (2) based on their medians and IQRs. You do not need to calculate these statistics; simply state how the medians and IQRs compare. Make sure to explain your reasoning.

- |   |  |
|---|--|
| (a) (1) 3, 5, 6, 7, 9<br>(2) 3, 5, 6, 7, 20 | (c) (1) 1, 2, 3, 4, 5<br>(2) 6, 7, 8, 9, 10              |
| (b) (1) 3, 5, 6, 7, 9<br>(2) 3, 5, 7, 8, 9  | (d) (1) 0, 10, 50, 60, 100<br>(2) 0, 100, 500, 600, 1000 |

**1.47 Means and SDs.** For each part, compare distributions (1) and (2) based on their means and standard deviations. You do not need to calculate these statistics; simply state how the means and the standard deviations compare. Make sure to explain your reasoning. *Hint:* It may be useful to sketch dot plots of the distributions.

- |  |   |
|--|---|
| (a) (1) 3, 5, 5, 5, 8, 11, 11, 11, 13<br>(2) 3, 5, 5, 5, 8, 11, 11, 11, 20 | (c) (1) 0, 2, 4, 6, 8, 10<br>(2) 20, 22, 24, 26, 28, 30     |
| (b) (1) -20, 0, 0, 0, 15, 25, 30, 30<br>(2) -40, 0, 0, 0, 15, 25, 30, 30   | (d) (1) 100, 200, 300, 400, 500<br>(2) 0, 50, 300, 550, 600 |

**1.48 Stats scores.** Below are the final exam scores of twenty introductory statistics students.

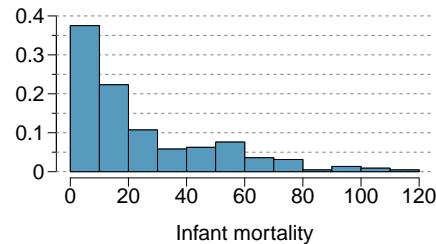
57, 66, 69, 71, 72, 73, 74, 77, 78, 78, 79, 79, 81, 81, 82, 83, 83, 88, 89, 94

Create a box plot of the distribution of these scores. The five number summary provided below may be useful.

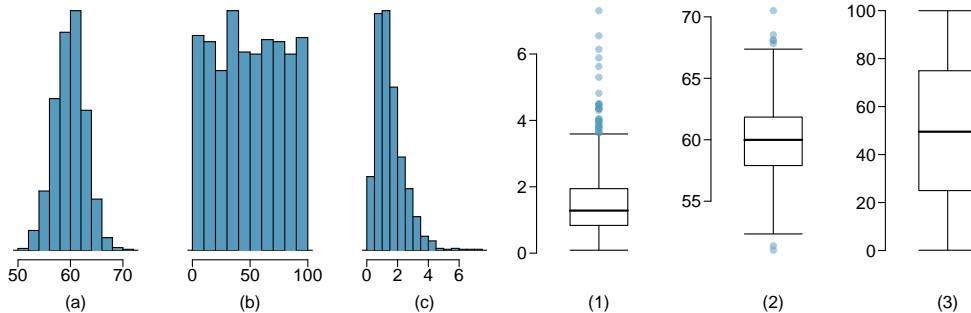
Min	Q1	Q2 (Median)	Q3	Max
57	72.5	78.5	82.5	94

**1.49 Infant mortality.** The infant mortality rate is defined as the number of infant deaths per 1,000 live births. This rate is often used as an indicator of the level of health in a country. The relative frequency histogram below shows the distribution of estimated infant death rates for 224 countries for which such data were available in 2014.<sup>67</sup>

- (a) Estimate Q1, the median, and Q3 from the histogram.
- (b) Would you expect the mean of this data set to be smaller or larger than the median? Explain your reasoning.

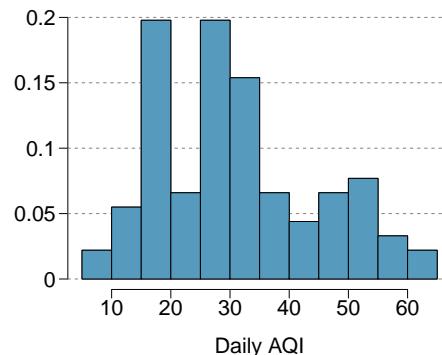


**1.50 Mix-and-match.** Describe the distribution in the histograms below and match them to the box plots.



**1.51 Air quality.** Daily air quality is measured by the air quality index (AQI) reported by the Environmental Protection Agency. This index reports the pollution level and what associated health effects might be a concern. The index is calculated for five major air pollutants regulated by the Clean Air Act and takes values from 0 to 300, where a higher value indicates lower air quality. AQI was reported for a sample of 91 days in 2011 in Durham, NC.<sup>68</sup> The relative frequency histogram below shows the distribution of the AQI values on these days.

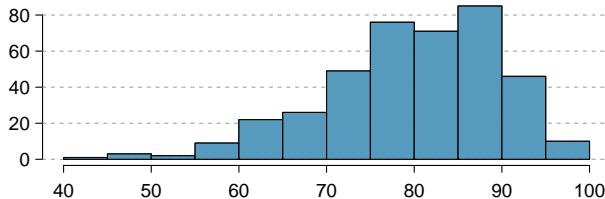
- (a) Estimate the median AQI value of this sample.
- (b) Would you expect the mean AQI value of this sample to be higher or lower than the median? Explain your reasoning.
- (c) Estimate Q1, Q3, and IQR for the distribution.
- (d) Would any of the days in this sample be considered to have an unusually low or high AQI? Explain your reasoning.



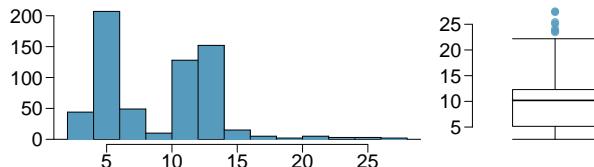
<sup>67</sup>CIA Factbook, Country Comparisons, 2014.

<sup>68</sup>US Environmental Protection Agency, AirData, 2011.

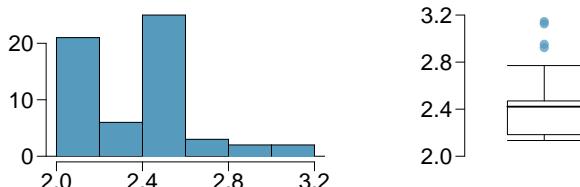
**1.52 Median vs. mean.** Estimate the median for the 400 observations shown in the histogram, and note whether you expect the mean to be higher or lower than the median.



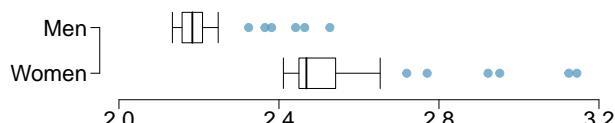
**1.53 Histograms vs. box plots.** Compare the two plots below. What characteristics of the distribution are apparent in the histogram and not in the box plot? What characteristics are apparent in the box plot but not in the histogram?



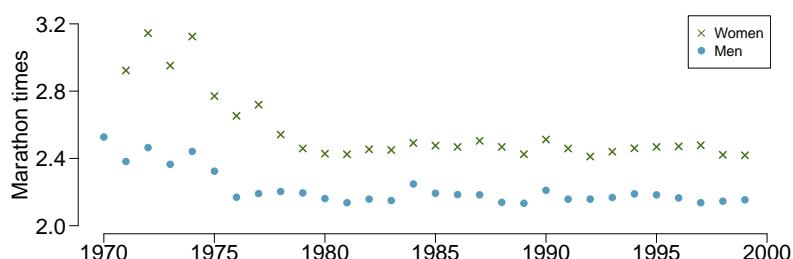
**1.54 Marathon winners.** The histogram and box plots below show the distribution of finishing times for male and female winners of the New York Marathon between 1970 and 1999.



- (a) What features of the distribution are apparent in the histogram and not the box plot? What features are apparent in the box plot but not in the histogram?
- (b) What may be the reason for the bimodal distribution? Explain.
- (c) Compare the distribution of marathon times for men and women based on the box plot shown below.



- (d) The time series plot shown below is another way to look at these data. Describe what is visible in this plot but not in the others.



**1.55 Distributions and appropriate statistics, Part I.** For each of the following, state whether you expect the distribution to be symmetric, right skewed, or left skewed. Also specify whether the mean or median would best represent a typical observation in the data, and whether the variability of observations would be best represented using the standard deviation or IQR. Explain your reasoning.

- (a) Number of pets per household.
- (b) Distance to work, i.e. number of miles between work and home.
- (c) Heights of adult males.

**1.56 Distributions and appropriate statistics, Part II.** For each of the following, state whether you expect the distribution to be symmetric, right skewed, or left skewed. Also specify whether the mean or median would best represent a typical observation in the data, and whether the variability of observations would be best represented using the standard deviation or IQR. Explain your reasoning.

- (a) Housing prices in a country where 25% of the houses cost below \$350,000, 50% of the houses cost below \$450,000, 75% of the houses cost below \$1,000,000 and there are a meaningful number of houses that cost more than \$6,000,000.
- (b) Housing prices in a country where 25% of the houses cost below \$300,000, 50% of the houses cost below \$600,000, 75% of the houses cost below \$900,000 and very few houses that cost more than \$1,200,000.
- (c) Number of alcoholic drinks consumed by college students in a given week. Assume that most of these students don't drink since they are under 21 years old, and only a few drink excessively.
- (d) Annual salaries of the employees at a Fortune 500 company where only a few high level executives earn much higher salaries than all the other employees.

**1.57 TV watchers.** Students in an AP Statistics class were asked how many hours of television they watch per week (including online streaming). This sample yielded an average of 4.71 hours, with a standard deviation of 4.18 hours. Is the distribution of number of hours students watch television weekly symmetric? If not, what shape would you expect this distribution to have? Explain your reasoning.

**1.58 Exam scores.** The average on a history exam (scored out of 100 points) was 85, with a standard deviation of 15. Is the distribution of the scores on this exam symmetric? If not, what shape would you expect this distribution to have? Explain your reasoning.

**1.59 Facebook friends.** Facebook data indicate that 50% of Facebook users have 100 or more friends, and that the average friend count of users is 190. What do these findings suggest about the shape of the distribution of number of friends of Facebook users?<sup>69</sup>

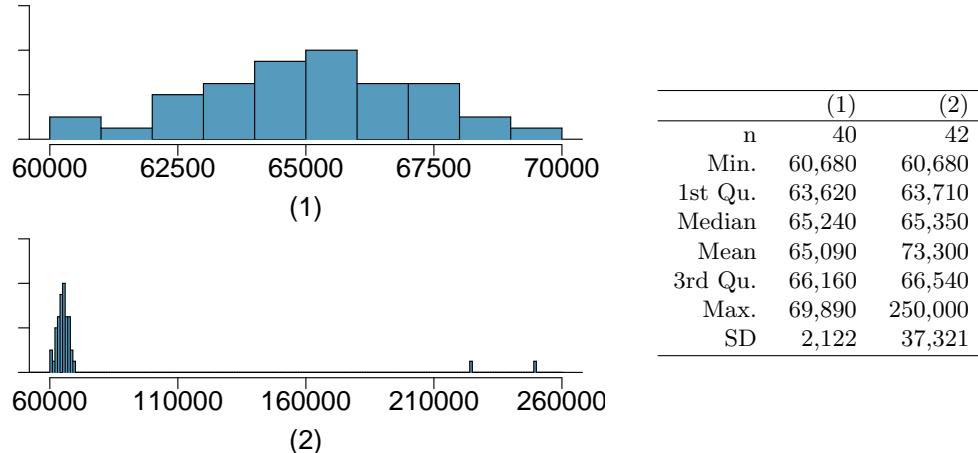
**1.60 A new statistic.** The statistic  $\frac{\bar{x}}{\text{median}}$  can be used as a measure of skewness. Suppose we have a distribution where all observations are greater than 0,  $x_i > 0$ . What is the expected shape of the distribution under the following conditions? Explain your reasoning.

- (a)  $\frac{\bar{x}}{\text{median}} = 1$
- (b)  $\frac{\bar{x}}{\text{median}} < 1$
- (c)  $\frac{\bar{x}}{\text{median}} > 1$

---

<sup>69</sup>Lars Backstrom. “Anatomy of Facebook”. In: *Facebook Data Teams Notes* (2011).

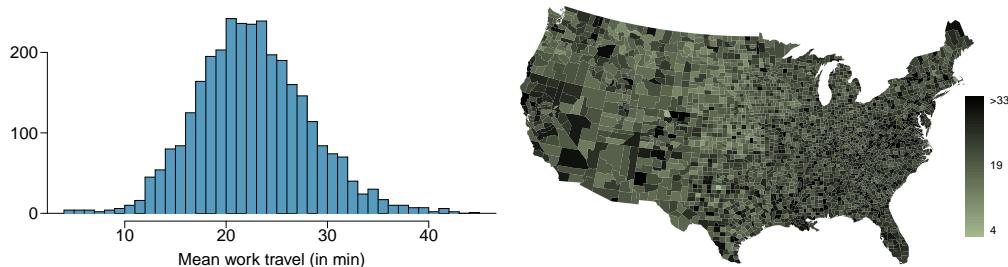
**1.61 Income at the coffee shop.** The first histogram below shows the distribution of the yearly incomes of 40 patrons at a college coffee shop. Suppose two new people walk into the coffee shop: one making \$225,000 and the other \$250,000. The second histogram shows the new income distribution. Summary statistics are also provided.



- Would the mean or the median best represent what we might think of as a typical income for the 42 patrons at this coffee shop? What does this say about the robustness of the two measures?
- Would the standard deviation or the IQR best represent the amount of variability in the incomes of the 42 patrons at this coffee shop? What does this say about the robustness of the two measures?

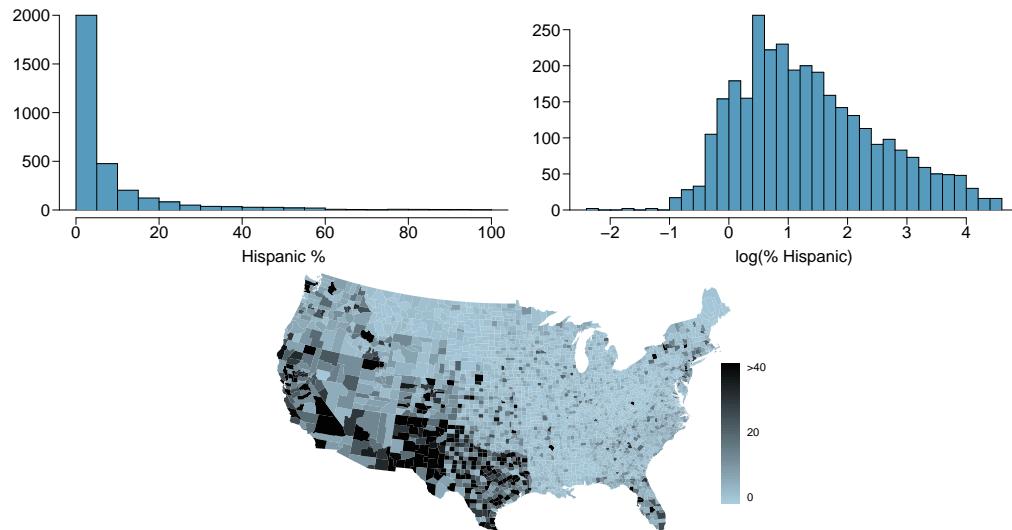
**1.62 Midrange.** The *midrange* of a distribution is defined as the average of the maximum and the minimum of that distribution. Is this statistic robust to outliers and extreme skew? Explain your reasoning

**1.63 Commute times.** The US census collects data on time it takes Americans to commute to work, among many other variables. The histogram below shows the distribution of average commute times in 3,143 US counties in 2010. Also shown below is a spatial intensity map of the same data.



- Describe the numerical distribution and comment on whether or not a log transformation may be advisable for these data.
- Describe the spatial distribution of commuting times using the map below.

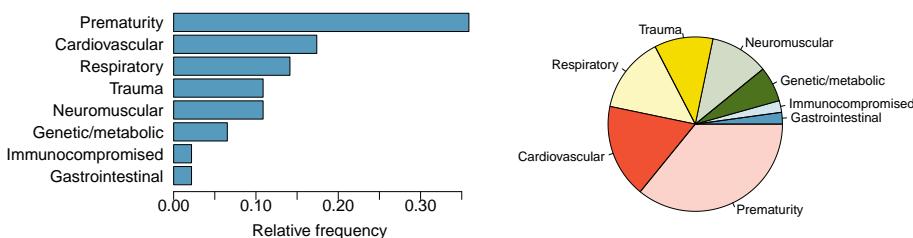
**1.64 Hispanic population.** The US census collects data on race and ethnicity of Americans, among many other variables. The histogram below shows the distribution of the percentage of the population that is Hispanic in 3,143 counties in the US in 2010. Also shown is a histogram of logs of these values.



- (a) Describe the numerical distribution and comment on why we might want to use log-transformed values in analyzing or modeling these data.
- (b) What features of the distribution of the Hispanic population in US counties are apparent in the map but not in the histogram? What features are apparent in the histogram but not the map?
- (c) Is one visualization more appropriate or helpful than the other? Explain your reasoning.

### 1.9.7 Considering categorical data

**1.65 Antibiotic use in children.** The bar plot and the pie chart below show the distribution of pre-existing medical conditions of children involved in a study on the optimal duration of antibiotic use in treatment of tracheitis, which is an upper respiratory infection.



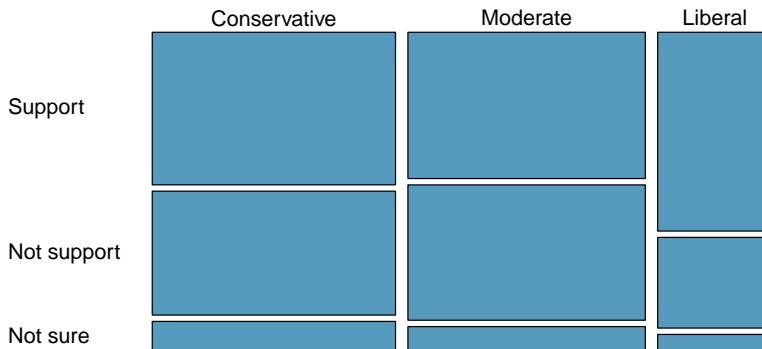
- (a) What features are apparent in the bar plot but not in the pie chart?
- (b) What features are apparent in the pie chart but not in the bar plot?
- (c) Which graph would you prefer to use for displaying these categorical data?

**1.66 Views on immigration.** 910 randomly sampled registered voters from Tampa, FL were asked if they thought workers who have illegally entered the US should be (i) allowed to keep their jobs and apply for US citizenship, (ii) allowed to keep their jobs as temporary guest workers but not allowed to apply for US citizenship, or (iii) lose their jobs and have to leave the country. The results of the survey by political ideology are shown below.<sup>70</sup>

		Political ideology			Total
		Conservative	Moderate	Liberal	
<i>Response</i>	(i) Apply for citizenship	57	120	101	278
	(ii) Guest worker	121	113	28	262
	(iii) Leave the country	179	126	45	350
	(iv) Not sure	15	4	1	20
	Total	372	363	175	910

- (a) What percent of these Tampa, FL voters identify themselves as conservatives?
- (b) What percent of these Tampa, FL voters are in favor of the citizenship option?
- (c) What percent of these Tampa, FL voters identify themselves as conservatives and are in favor of the citizenship option?
- (d) What percent of these Tampa, FL voters who identify themselves as conservatives are also in favor of the citizenship option? What percent of moderates share this view? What percent of liberals share this view?
- (e) Do political ideology and views on immigration appear to be independent? Explain your reasoning.

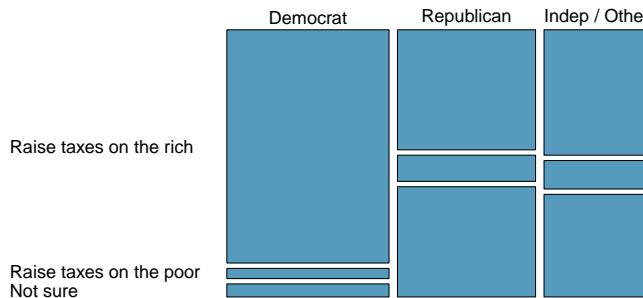
**1.67 Views on the DREAM Act.** A random sample of registered voters from Tampa, FL were asked if they support the DREAM Act, a proposed law which would provide a path to citizenship for people brought illegally to the US as children. The survey also collected information on the political ideology of the respondents. Based on the mosaic plot shown below, do views on the DREAM Act and political ideology appear to be independent? Explain your reasoning.<sup>71</sup>



<sup>70</sup>SurveyUSA, News Poll #18927, data collected Jan 27-29, 2012.

<sup>71</sup>SurveyUSA, News Poll #18927, data collected Jan 27-29, 2012.

**1.68 Raise taxes.** A random sample of registered voters nationally were asked whether they think it's better to raise taxes on the rich or raise taxes on the poor. The survey also collected information on the political party affiliation of the respondents. Based on the mosaic plot shown below, do views on raising taxes and political affiliation appear to be independent? Explain your reasoning.<sup>72</sup>



### 1.9.8 Case study: gender discrimination

**1.69 Side effects of Avandia.** Rosiglitazone is the active ingredient in the controversial type 2 diabetes medicine Avandia and has been linked to an increased risk of serious cardiovascular problems such as stroke, heart failure, and death. A common alternative treatment is pioglitazone, the active ingredient in a diabetes medicine called Actos. In a nationwide retrospective observational study of 227,571 Medicare beneficiaries aged 65 years or older, it was found that 2,593 of the 67,593 patients using rosiglitazone and 5,386 of the 159,978 using pioglitazone had serious cardiovascular problems. These data are summarized in the contingency table below.<sup>73</sup>

		<i>Cardiovascular problems</i>		Total
		Yes	No	
Treatment	Rosiglitazone	2,593	65,000	67,593
	Pioglitazone	5,386	154,592	159,978
	Total	7,979	219,592	227,571

- (a) Determine if each of the following statements is true or false. If false, explain why. *Be careful:* The reasoning may be wrong even if the statement's conclusion is correct. In such cases, the statement should be considered false.
- i. Since more patients on pioglitazone had cardiovascular problems (5,386 vs. 2,593), we can conclude that the rate of cardiovascular problems for those on a pioglitazone treatment is higher.
  - ii. The data suggest that diabetic patients who are taking rosiglitazone are more likely to have cardiovascular problems since the rate of incidence was ( $2,593 / 67,593 = 0.038$ ) 3.8% for patients on this treatment, while it was only ( $5,386 / 159,978 = 0.034$ ) 3.4% for patients on pioglitazone.
  - iii. The fact that the rate of incidence is higher for the rosiglitazone group proves that rosiglitazone causes serious cardiovascular problems.
  - iv. Based on the information provided so far, we cannot tell if the difference between the rates of incidences is due to a relationship between the two variables or due to chance.

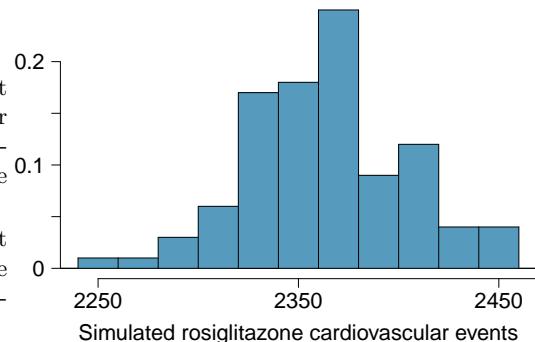
**(See the next page for additional parts to this question.)**

<sup>72</sup>Public Policy Polling, Americans on College Degrees, Classic Literature, the Seasons, and More, data collected Feb 20-22, 2015.

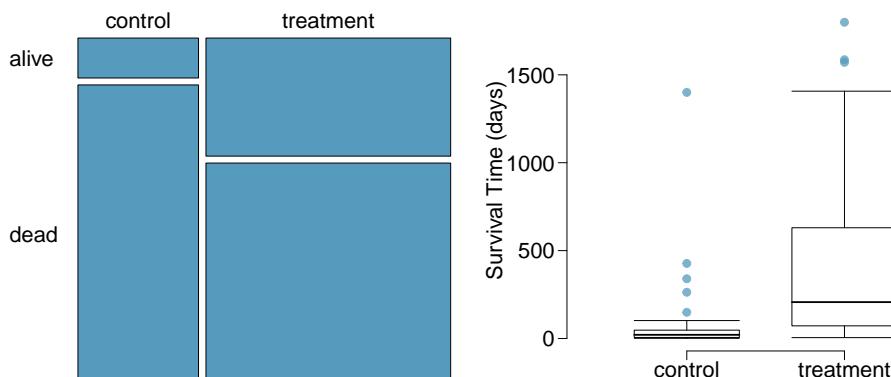
<sup>73</sup>D.J. Graham et al. "Risk of acute myocardial infarction, stroke, heart failure, and death in elderly Medicare patients treated with rosiglitazone or pioglitazone". In: *JAMA* 304.4 (2010), p. 411. ISSN: 0098-7484.

- (b) What proportion of all patients had cardiovascular problems?
- (c) If the type of treatment and having cardiovascular problems were independent, about how many patients in the rosiglitazone group would we expect to have had cardiovascular problems?
- (d) We can investigate the relationship between outcome and treatment in this study using a randomization technique. While in reality we would carry out the simulations required for randomization using statistical software, suppose we actually simulate using index cards. In order to simulate from the independence model, which states that the outcomes were independent of the treatment, we write whether or not each patient had a cardiovascular problem on cards, shuffled all the cards together, then deal them into two groups of size 67,593 and 159,978. We repeat this simulation 1,000 times and each time record the number of people in the rosiglitazone group who had cardiovascular problems. Use the relative frequency histogram of these counts to answer (i)-(iii).

- i. What are the claims being tested?
- ii. Compared to the number calculated in part (b), which would provide more support for the alternative hypothesis, *more* or *fewer* patients with cardiovascular problems in the rosiglitazone group?
- iii. What do the simulation results suggest about the relationship between taking rosiglitazone and having cardiovascular problems in diabetic patients?



**1.70 Heart transplants.** The Stanford University Heart Transplant Study was conducted to determine whether an experimental heart transplant program increased lifespan. Each patient entering the program was designated an official heart transplant candidate, meaning that he was gravely ill and would most likely benefit from a new heart. Some patients got a transplant and some did not. The variable `transplant` indicates which group the patients were in; patients in the treatment group got a transplant and those in the control group did not. Another variable called `survived` was used to indicate whether or not the patient was alive at the end of the study. Of the 34 patients in the control group, 30 died. Of the 69 people in the treatment group, 45 died.<sup>74</sup>



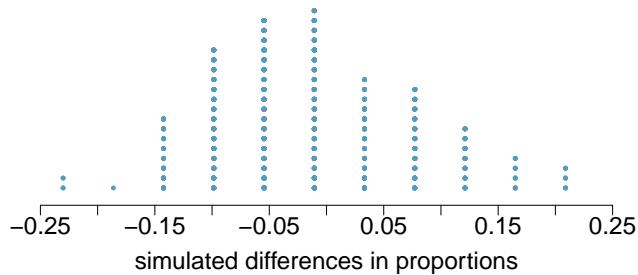
- (a) Based on the mosaic plot, is survival independent of whether or not the patient got a transplant? Explain your reasoning.  
**(See the next page for additional parts to this question.)**

<sup>74</sup>B. Turnbull et al. "Survivorship of Heart Transplant Data". In: *Journal of the American Statistical Association* 69 (1974), pp. 74–80.

- (b) What do the box plots below suggest about the efficacy (effectiveness) of the heart transplant treatment.
- (c) What proportion of patients in the treatment group and what proportion of patients in the control group died?
- (d) One approach for investigating whether or not the treatment is effective is to use a randomization technique.
- What are the claims being tested?
  - The paragraph below describes the set up for such approach, if we were to do it without using statistical software. Fill in the blanks with a number or phrase, whichever is appropriate.

We write *alive* on \_\_\_\_\_ cards representing patients who were alive at the end of the study, and *dead* on \_\_\_\_\_ cards representing patients who were not. Then, we shuffle these cards and split them into two groups: one group of size \_\_\_\_\_ representing treatment, and another group of size \_\_\_\_\_ representing control. We calculate the difference between the proportion of *dead* cards in the treatment and control groups (treatment - control) and record this value. We repeat this 100 times to build a distribution centered at \_\_\_\_\_. Lastly, we calculate the fraction of simulations where the simulated differences in proportions are \_\_\_\_\_. If this fraction is low, we conclude that it is unlikely to have observed such an outcome by chance and that the null hypothesis should be rejected in favor of the alternative.

- What do the simulation results shown below suggest about the effectiveness of the transplant program?



## Chapter 3

# Distributions of random variables

### 3.1 Normal distribution

Among all the distributions we see in practice, one is overwhelmingly the most common. The symmetric, unimodal, bell curve is ubiquitous throughout statistics. Indeed it is so common, that people often know it as the **normal curve** or **normal distribution**,<sup>1</sup> shown in Figure 3.1. Variables such as SAT scores and heights of US adult males closely follow the normal distribution.

#### Normal distribution facts

Many variables are nearly normal, but none are exactly normal. Thus the normal distribution, while not perfect for any single problem, is very useful for a variety of problems. We will use it in data exploration and to solve important problems in statistics.

---

<sup>1</sup>It is also introduced as the Gaussian distribution after Frederic Gauss, the first person to formalize its mathematical expression.

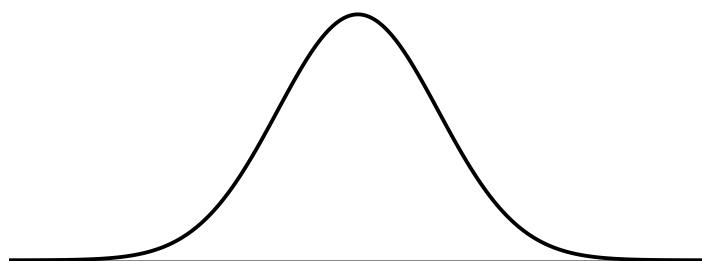


Figure 3.1: A normal curve.

### 3.1.1 Normal distribution model

The normal distribution model always describes a symmetric, unimodal, bell-shaped curve. However, these curves can look different depending on the details of the model. Specifically, the normal distribution model can be adjusted using two parameters: mean and standard deviation. As you can probably guess, changing the mean shifts the bell curve to the left or right, while changing the standard deviation stretches or constricts the curve. Figure 3.2 shows the normal distribution with mean 0 and standard deviation 1 in the left panel and the normal distributions with mean 19 and standard deviation 4 in the right panel. Figure 3.3 shows these distributions on the same axis.

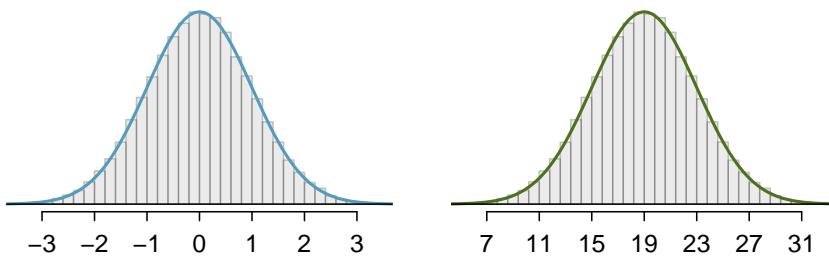


Figure 3.2: Both curves represent the normal distribution, however, they differ in their center and spread. The normal distribution with mean 0 and standard deviation 1 is called the **standard normal distribution**.

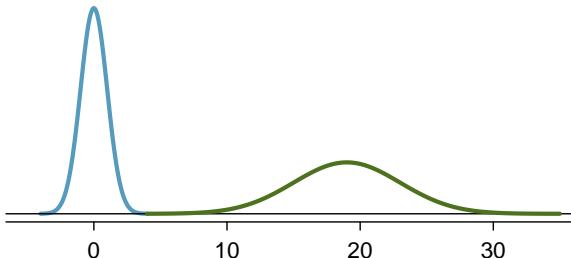


Figure 3.3: The normal models shown in Figure 3.2 but plotted together and on the same scale.

If a normal distribution has mean  $\mu$  and standard deviation  $\sigma$ , we may write the distribution as  $N(\mu, \sigma)$ . The two distributions in Figure 3.3 can be written as

$$N(\mu = 0, \sigma = 1) \quad \text{and} \quad N(\mu = 19, \sigma = 4)$$

Because the mean and standard deviation describe a normal distribution exactly, they are called the distribution's **parameters**.

 **Guided Practice 3.1** Write down the short-hand for a normal distribution with<sup>2</sup>

- (a) mean 5 and standard deviation 3,
- (b) mean -100 and standard deviation 10, and
- (c) mean 2 and standard deviation 9.

$N(\mu, \sigma)$   
Normal dist.  
with mean  $\mu$   
& st. dev.  $\sigma$

---

<sup>2</sup>(a)  $N(\mu = 5, \sigma = 3)$ . (b)  $N(\mu = -100, \sigma = 10)$ . (c)  $N(\mu = 2, \sigma = 9)$ .

	SAT	ACT
Mean	1500	21
SD	300	5

Table 3.4: Mean and standard deviation for the SAT and ACT.

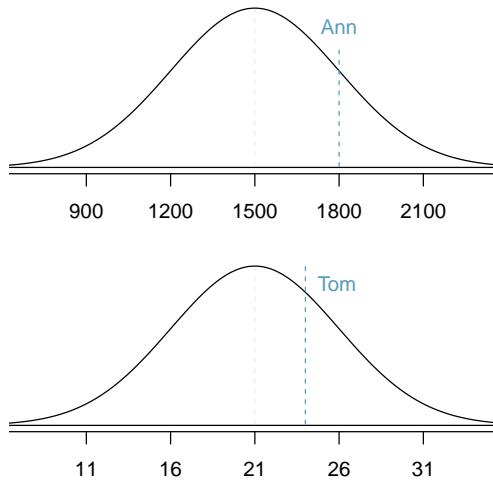


Figure 3.5: Ann's and Tom's scores shown with the distributions of SAT and ACT scores.

### 3.1.2 Standardizing with Z-scores

- Example 3.2 Table 3.4 shows the mean and standard deviation for total scores on the SAT and ACT. The distribution of SAT and ACT scores are both nearly normal. Suppose Ann scored 1800 on her SAT and Tom scored 24 on his ACT. Who performed better?

We use the standard deviation as a guide. Ann is 1 standard deviation above average on the SAT:  $1500 + 300 = 1800$ . Tom is 0.6 standard deviations above the mean on the ACT:  $21 + 0.6 \times 5 = 24$ . In Figure 3.5, we can see that Ann tends to do better with respect to everyone else than Tom did, so her score was better.

Example 3.2 used a standardization technique called a Z-score, a method most commonly employed for nearly normal observations but that may be used with any distribution. The **Z-score** of an observation is defined as the number of standard deviations it falls above or below the mean. If the observation is one standard deviation above the mean, its Z-score is 1. If it is 1.5 standard deviations *below* the mean, then its Z-score is -1.5. If  $x$  is an observation from a distribution  $N(\mu, \sigma)$ , we define the Z-score mathematically as

$$Z = \frac{x - \mu}{\sigma}$$

Using  $\mu_{SAT} = 1500$ ,  $\sigma_{SAT} = 300$ , and  $x_{Ann} = 1800$ , we find Ann's Z-score:

$$Z_{Ann} = \frac{x_{Ann} - \mu_{SAT}}{\sigma_{SAT}} = \frac{1800 - 1500}{300} = 1$$

**Z**  
Z-score, the  
standardized  
observation

### The Z-score

The Z-score of an observation is the number of standard deviations it falls above or below the mean. We compute the Z-score for an observation  $x$  that follows a distribution with mean  $\mu$  and standard deviation  $\sigma$  using

$$Z = \frac{x - \mu}{\sigma}$$

- **Guided Practice 3.3** Use Tom's ACT score, 24, along with the ACT mean and standard deviation to compute his Z-score.<sup>3</sup>

Observations above the mean always have positive Z-scores while those below the mean have negative Z-scores. If an observation is equal to the mean (e.g. SAT score of 1500), then the Z-score is 0.

- **Guided Practice 3.4** Let  $X$  represent a random variable from  $N(\mu = 3, \sigma = 2)$ , and suppose we observe  $x = 5.19$ . (a) Find the Z-score of  $x$ . (b) Use the Z-score to determine how many standard deviations above or below the mean  $x$  falls.<sup>4</sup>

- **Guided Practice 3.5** Head lengths of brushtail possums follow a nearly normal distribution with mean 92.6 mm and standard deviation 3.6 mm. Compute the Z-scores for possums with head lengths of 95.4 mm and 85.8 mm.<sup>5</sup>

We can use Z-scores to roughly identify which observations are more unusual than others. One observation  $x_1$  is said to be more unusual than another observation  $x_2$  if the absolute value of its Z-score is larger than the absolute value of the other observation's Z-score:  $|Z_1| > |Z_2|$ . This technique is especially insightful when a distribution is symmetric.

- **Guided Practice 3.6** Which of the observations in Guided Practice 3.5 is more unusual?<sup>6</sup>

### 3.1.3 Normal probability table

- **Example 3.7** Ann from Example 3.2 earned a score of 1800 on her SAT with a corresponding  $Z = 1$ . She would like to know what percentile she falls in among all SAT test-takers.

Ann's **percentile** is the percentage of people who earned a lower SAT score than Ann. We shade the area representing those individuals in Figure 3.6. The total area under the normal curve is always equal to 1, and the proportion of people who scored below Ann on the SAT is equal to the *area* shaded in Figure 3.6: 0.8413. In other words, Ann is in the 84<sup>th</sup> percentile of SAT takers.

---

<sup>3</sup> $Z_{Tom} = \frac{x_{Tom} - \mu_{ACT}}{\sigma_{ACT}} = \frac{24 - 21}{5} = 0.6$

<sup>4</sup>(a) Its Z-score is given by  $Z = \frac{x - \mu}{\sigma} = \frac{5.19 - 3}{2} = 2.19/2 = 1.095$ . (b) The observation  $x$  is 1.095 standard deviations *above* the mean. We know it must be above the mean since  $Z$  is positive.

<sup>5</sup>For  $x_1 = 95.4$  mm:  $Z_1 = \frac{x_1 - \mu}{\sigma} = \frac{95.4 - 92.6}{3.6} = 0.78$ . For  $x_2 = 85.8$  mm:  $Z_2 = \frac{85.8 - 92.6}{3.6} = -1.89$ .

<sup>6</sup>Because the *absolute value* of Z-score for the second observation is larger than that of the first, the second observation has a more unusual head length.

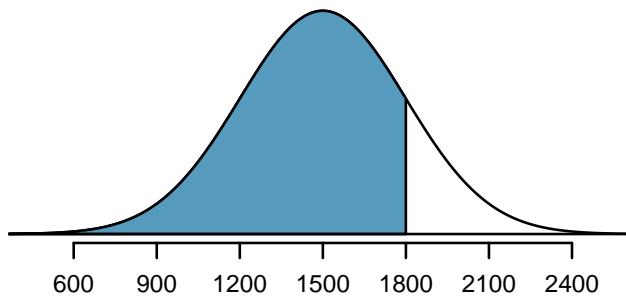


Figure 3.6: The normal model for SAT scores, shading the area of those individuals who scored below Ann.

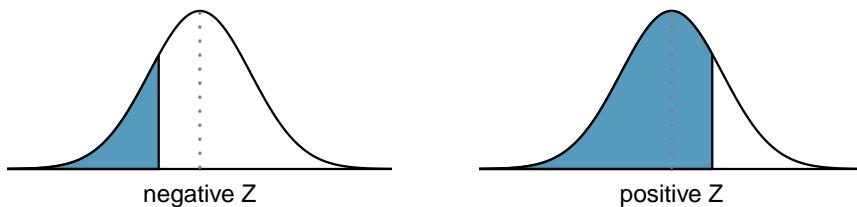


Figure 3.7: The area to the left of  $Z$  represents the percentile of the observation.

We can use the normal model to find percentiles. A **normal probability table**, which lists  $Z$ -scores and corresponding percentiles, can be used to identify a percentile based on the  $Z$ -score (and vice versa). Statistical software can also be used.

A normal probability table is given in Appendix B.1 on page 252 and abbreviated in Table 3.8. We use this table to identify the percentile corresponding to any particular  $Z$ -score. For instance, the percentile of  $Z = 0.43$  is shown in row 0.4 and column 0.03 in Table 3.8: 0.6664, or the 66.64<sup>th</sup> percentile. Generally, we round  $Z$  to two decimals, identify the proper row in the normal probability table up through the first decimal, and then determine the column representing the second decimal value. The intersection of this row and column is the percentile of the observation.

We can also find the  $Z$ -score associated with a percentile. For example, to identify  $Z$  for the 80<sup>th</sup> percentile, we look for the value closest to 0.8000 in the middle portion of the table: 0.7995. We determine the  $Z$ -score for the 80<sup>th</sup> percentile by combining the row and column  $Z$  values: 0.84.

- **Guided Practice 3.8** Determine the proportion of SAT test takers who scored better than Ann on the SAT.<sup>7</sup>

---

<sup>7</sup>If 84% had lower scores than Ann, the proportion of people who had better scores must be 16%. (Generally ties are ignored when the normal model, or any other continuous distribution, is used.)

Z	Second decimal place of Z									
	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

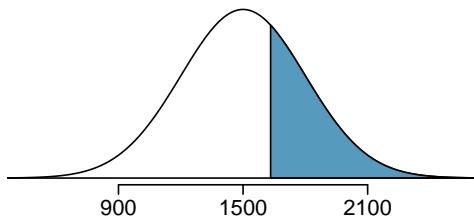
Table 3.8: A section of the normal probability table. The percentile for a normal random variable with  $Z = 0.43$  has been *highlighted*, and the percentile closest to 0.8000 has also been *highlighted*.

### 3.1.4 Normal probability examples

Cumulative SAT scores are approximated well by a normal model,  $N(\mu = 1500, \sigma = 300)$ .

- Example 3.9 Shannon is a randomly selected SAT taker, and nothing is known about Shannon's SAT aptitude. What is the probability Shannon scores at least 1630 on her SATs?

First, always draw and label a picture of the normal distribution. (Drawings need not be exact to be useful.) We are interested in the chance she scores above 1630, so we shade this upper tail:



The picture shows the mean and the values at 2 standard deviations above and below the mean. The simplest way to find the shaded area under the curve makes use of the Z-score of the cutoff value. With  $\mu = 1500$ ,  $\sigma = 300$ , and the cutoff value  $x = 1630$ , the Z-score is computed as

$$Z = \frac{x - \mu}{\sigma} = \frac{1630 - 1500}{300} = \frac{130}{300} = 0.43$$

We look up the percentile of  $Z = 0.43$  in the normal probability table shown in Table 3.8 or in Appendix B.1 on page 252, which yields 0.6664. However, the percentile

describes those who had a Z-score *lower* than 0.43. To find the area *above*  $Z = 0.43$ , we compute one minus the area of the lower tail:

$$1.0000 - 0.6664 = 0.3336$$

The probability Shannon scores at least 1630 on the SAT is 0.3336.

**TIP: always draw a picture first, and find the Z-score second**

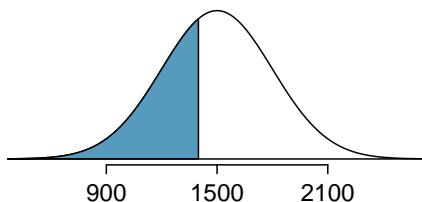
For any normal probability situation, *always always always* draw and label the normal curve and shade the area of interest first. The picture will provide an estimate of the probability.

After drawing a figure to represent the situation, identify the Z-score for the observation of interest.

- ④ **Guided Practice 3.10** If the probability of Shannon scoring at least 1630 is 0.3336, then what is the probability she scores less than 1630? Draw the normal curve representing this exercise, shading the lower region instead of the upper one.<sup>8</sup>

- ⑤ **Example 3.11** Edward earned a 1400 on his SAT. What is his percentile?

First, a picture is needed. Edward's percentile is the proportion of people who do not get as high as a 1400. These are the scores to the left of 1400.



Identifying the mean  $\mu = 1500$ , the standard deviation  $\sigma = 300$ , and the cutoff for the tail area  $x = 1400$  makes it easy to compute the Z-score:

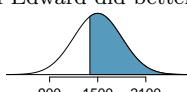
$$Z = \frac{x - \mu}{\sigma} = \frac{1400 - 1500}{300} = -0.33$$

Using the normal probability table, identify the row of  $-0.3$  and column of  $0.03$ , which corresponds to the probability 0.3707. Edward is at the 37<sup>th</sup> percentile.

- ⑥ **Guided Practice 3.12** Use the results of Example 3.11 to compute the proportion of SAT takers who did better than Edward. Also draw a new picture.<sup>9</sup>

<sup>8</sup>We found the probability in Example 3.9: 0.6664. A picture for this exercise is represented by the shaded area below "0.6664" in Example 3.9.

<sup>9</sup>If Edward did better than 37% of SAT takers, then about 63% must have done better than him.



**TIP: areas to the right**

The normal probability table in most books gives the area to the left. If you would like the area to the right, first find the area to the left and then subtract this amount from one.

- Ⓐ **Guided Practice 3.13** Stuart earned an SAT score of 2100. Draw a picture for each part. (a) What is his percentile? (b) What percent of SAT takers did better than Stuart?<sup>10</sup>

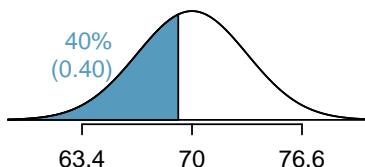
Based on a sample of 100 men,<sup>11</sup> the heights of male adults between the ages 20 and 62 in the US is nearly normal with mean 70.0" and standard deviation 3.3".

- Ⓑ **Guided Practice 3.14** Mike is 5'7" and Jim is 6'4". (a) What is Mike's height percentile? (b) What is Jim's height percentile? Also draw one picture for each part.<sup>12</sup>

The last several problems have focused on finding the probability or percentile for a particular observation. What if you would like to know the observation corresponding to a particular percentile?

- **Example 3.15** Erik's height is at the 40<sup>th</sup> percentile. How tall is he?

As always, first draw the picture.



In this case, the lower tail probability is known (0.40), which can be shaded on the diagram. We want to find the observation that corresponds to this value. As a first step in this direction, we determine the Z-score associated with the 40<sup>th</sup> percentile.

Because the percentile is below 50%, we know  $Z$  will be negative. Looking in the negative part of the normal probability table, we search for the probability *inside* the table closest to 0.4000. We find that 0.4000 falls in row -0.2 and between columns 0.05 and 0.06. Since it falls closer to 0.05, we take this one:  $Z = -0.25$ .

Knowing  $Z_{Erik} = -0.25$  and the population parameters  $\mu = 70$  and  $\sigma = 3.3$  inches, the Z-score formula can be set up to determine Erik's unknown height, labeled  $x_{Erik}$ :

$$-0.25 = Z_{Erik} = \frac{x_{Erik} - \mu}{\sigma} = \frac{x_{Erik} - 70}{3.3}$$

Solving for  $x_{Erik}$  yields the height 69.18 inches. That is, Erik is about 5'9" (this is notation for 5-feet, 9-inches).

<sup>10</sup>Numerical answers: (a) 0.9772. (b) 0.0228.

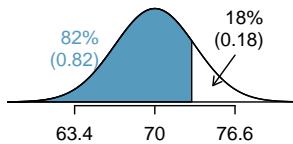
<sup>11</sup>This sample was taken from the USDA Food Commodity Intake Database.

<sup>12</sup>First put the heights into inches: 67 and 76 inches. Figures are shown below. (a)  $Z_{Mike} = \frac{67-70}{3.3} = -0.91 \rightarrow 0.1814$ . (b)  $Z_{Jim} = \frac{76-70}{3.3} = 1.82 \rightarrow 0.9656$ .



● **Example 3.16** What is the adult male height at the 82<sup>nd</sup> percentile?

Again, we draw the figure first.



Next, we want to find the Z-score at the 82<sup>nd</sup> percentile, which will be a positive value. Looking in the Z-table, we find  $Z$  falls in row 0.9 and the nearest column is 0.02, i.e.  $Z = 0.92$ . Finally, the height  $x$  is found using the Z-score formula with the known mean  $\mu$ , standard deviation  $\sigma$ , and Z-score  $Z = 0.92$ :

$$0.92 = Z = \frac{x - \mu}{\sigma} = \frac{x - 70}{3.3}$$

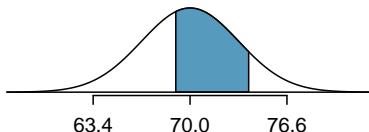
This yields 73.04 inches or about 6'1" as the height at the 82<sup>nd</sup> percentile.

○ **Guided Practice 3.17** (a) What is the 95<sup>th</sup> percentile for SAT scores? (b) What is the 97.5<sup>th</sup> percentile of the male heights? As always with normal probability problems, first draw a picture.<sup>13</sup>

○ **Guided Practice 3.18** (a) What is the probability that a randomly selected male adult is at least 6'2" (74 inches)? (b) What is the probability that a male adult is shorter than 5'9" (69 inches)?<sup>14</sup>

● **Example 3.19** What is the probability that a random adult male is between 5'9" and 6'2"?

These heights correspond to 69 inches and 74 inches. First, draw the figure. The area of interest is no longer an upper or lower tail.



The total area under the curve is 1. If we find the area of the two tails that are not shaded (from Guided Practice 3.18, these areas are 0.3821 and 0.1131), then we can find the middle area:

$$1.0000 - 0.3821 - 0.1131 = 0.5048$$



That is, the probability of being between 5'9" and 6'2" is 0.5048.

<sup>13</sup>Remember: draw a picture first, then find the Z-score. (We leave the pictures to you.) The Z-score can be found by using the percentiles and the normal probability table. (a) We look for 0.95 in the probability portion (middle part) of the normal probability table, which leads us to row 1.6 and (about) column 0.05, i.e.  $Z_{95} = 1.65$ . Knowing  $Z_{95} = 1.65$ ,  $\mu = 1500$ , and  $\sigma = 300$ , we setup the Z-score formula:  $1.65 = \frac{x_{95} - 1500}{300}$ . We solve for  $x_{95}$ :  $x_{95} = 1995$ . (b) Similarly, we find  $Z_{97.5} = 1.96$ , again setup the Z-score formula for the heights, and calculate  $x_{97.5} = 76.5$ .

<sup>14</sup>Numerical answers: (a) 0.1131. (b) 0.3821.

Ⓐ **Guided Practice 3.20** What percent of SAT takers get between 1500 and 2000?<sup>15</sup>

Ⓐ **Guided Practice 3.21** What percent of adult males are between 5'5" and 5'7"?<sup>16</sup>

### Calculator videos

Videos covering calculations for the normal distribution using TI and Casio graphing calculators are available at [openintro.org/videos](http://openintro.org/videos).

#### 3.1.5 68-95-99.7 rule

Here, we present a useful rule of thumb for the probability of falling within 1, 2, and 3 standard deviations of the mean in the normal distribution. This will be useful in a wide range of practical settings, especially when trying to make a quick estimate without a calculator or Z-table.

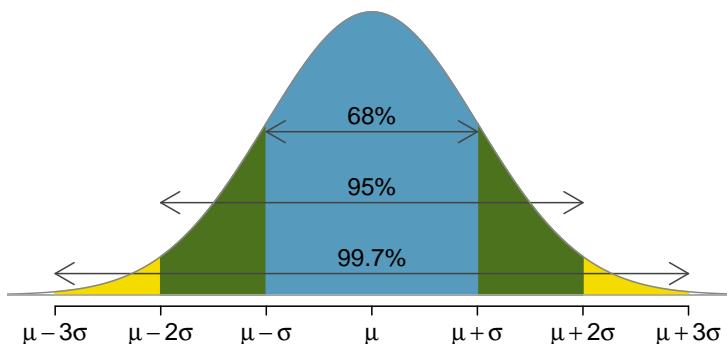


Figure 3.9: Probabilities for falling within 1, 2, and 3 standard deviations of the mean in a normal distribution.

Ⓐ **Guided Practice 3.22** Use the Z-table to confirm that about 68%, 95%, and 99.7% of observations fall within 1, 2, and 3, standard deviations of the mean in the normal distribution, respectively. For instance, first find the area that falls between  $Z = -1$  and  $Z = 1$ , which should have an area of about 0.68. Similarly there should be an area of about 0.95 between  $Z = -2$  and  $Z = 2$ .<sup>17</sup>

It is possible for a normal random variable to fall 4, 5, or even more standard deviations from the mean. However, these occurrences are very rare if the data are nearly normal. The probability of being further than 4 standard deviations from the mean is about 1-in-15,000. For 5 and 6 standard deviations, it is about 1-in-2 million and 1-in-500 million, respectively.

<sup>15</sup>This is an abbreviated solution. (Be sure to draw a figure!) First find the percent who get below 1500 and the percent that get above 2000:  $Z_{1500} = 0.00 \rightarrow 0.5000$  (area below),  $Z_{2000} = 1.67 \rightarrow 0.0475$  (area above). Final answer:  $1.0000 - 0.5000 - 0.0475 = 0.4525$ .

<sup>16</sup>5'5" is 65 inches. 5'7" is 67 inches. Numerical solution:  $1.000 - 0.0649 - 0.8183 = 0.1168$ , i.e. 11.68%.

<sup>17</sup>First draw the pictures. To find the area between  $Z = -1$  and  $Z = 1$ , use the normal probability table to determine the areas below  $Z = -1$  and above  $Z = 1$ . Next verify the area between  $Z = -1$  and  $Z = 1$  is about 0.68. Repeat this for  $Z = -2$  to  $Z = 2$  and also for  $Z = -3$  to  $Z = 3$ .

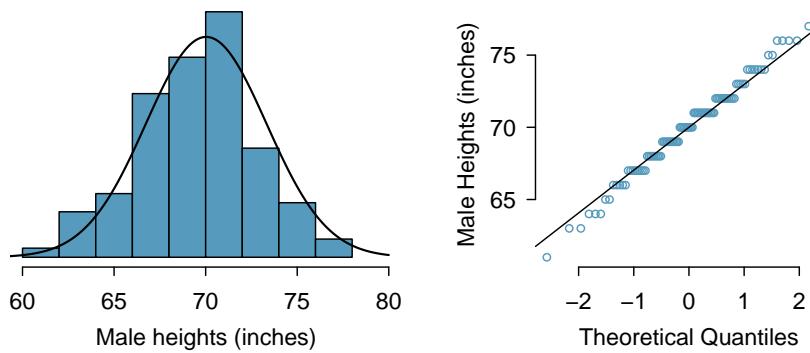


Figure 3.10: A sample of 100 male heights. The observations are rounded to the nearest whole inch, explaining why the points appear to jump in increments in the normal probability plot.

- Ⓐ **Guided Practice 3.23** SAT scores closely follow the normal model with mean  $\mu = 1500$  and standard deviation  $\sigma = 300$ . (a) About what percent of test takers score 900 to 2100? (b) What percent score between 1500 and 2100?<sup>18</sup>

## 3.2 Evaluating the normal approximation



Many processes can be well approximated by the normal distribution. We have already seen two good examples: SAT scores and the heights of US adult males. While using a normal model can be extremely convenient and helpful, it is important to remember normality is always an approximation. Evaluating the appropriateness of the normal assumption is a key step in many data analyses.

Example 3.15 suggests the distribution of heights of US males is well approximated by the normal model. We are interested in proceeding under the assumption that the data are normally distributed, but first we must check to see if this is reasonable.

There are two visual methods for checking the assumption of normality, which can be implemented and interpreted quickly. The first is a simple histogram with the best fitting normal curve overlaid on the plot, as shown in the left panel of Figure 3.10. The sample mean  $\bar{x}$  and standard deviation  $s$  are used as the parameters of the best fitting normal curve. The closer this curve fits the histogram, the more reasonable the normal model assumption. Another more common method is examining a **normal probability plot**,<sup>19</sup> shown in the right panel of Figure 3.10. The closer the points are to a perfect straight line, the more confident we can be that the data follow the normal model.

<sup>18</sup>(a) 900 and 2100 represent two standard deviations above and below the mean, which means about 95% of test takers will score between 900 and 2100. (b) Since the normal model is symmetric, then half of the test takers from part (a) ( $\frac{95\%}{2} = 47.5\%$  of all test takers) will score 900 to 1500 while 47.5% score between 1500 and 2100.

<sup>19</sup>Also commonly called a **quantile-quantile plot**.

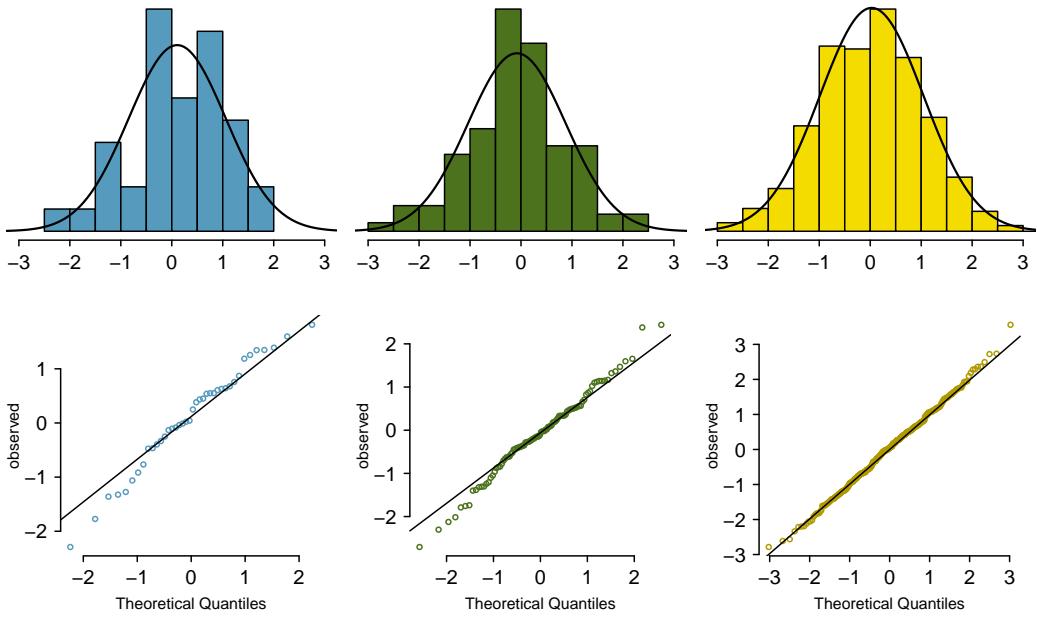


Figure 3.11: Histograms and normal probability plots for three simulated normal data sets;  $n = 40$  (left),  $n = 100$  (middle),  $n = 400$  (right).

**Example 3.24** Three data sets of 40, 100, and 400 samples were simulated from a normal distribution, and the histograms and normal probability plots of the data sets are shown in Figure 3.11. These will provide a benchmark for what to look for in plots of real data.

---

The left panels show the histogram (top) and normal probability plot (bottom) for the simulated data set with 40 observations. The data set is too small to really see clear structure in the histogram. The normal probability plot also reflects this, where there are some deviations from the line. We should expect deviations of this amount for such a small data set.

The middle panels show diagnostic plots for the data set with 100 simulated observations. The histogram shows more normality and the normal probability plot shows a better fit. While there are a few observations that deviate noticeably from the line, they are not particularly extreme.

The data set with 400 observations has a histogram that greatly resembles the normal distribution, while the normal probability plot is nearly a perfect straight line. Again in the normal probability plot there is one observation (the largest) that deviates slightly from the line. If that observation had deviated 3 times further from the line, it would be of greater importance in a real data set. Apparent outliers can occur in normally distributed data but they are rare.

Notice the histograms look more normal as the sample size increases, and the normal probability plot becomes straighter and more stable.

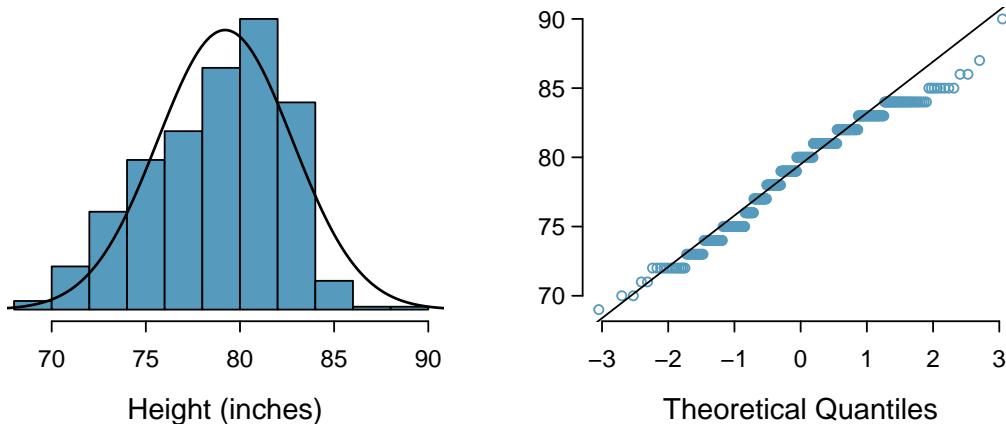


Figure 3.12: Histogram and normal probability plot for the NBA heights from the 2008-9 season.

- **Example 3.25** Are NBA player heights normally distributed? Consider all 435 NBA players from the 2008-9 season presented in Figure 3.12.<sup>20</sup>

We first create a histogram and normal probability plot of the NBA player heights. The histogram in the left panel is slightly left skewed, which contrasts with the symmetric normal distribution. The points in the normal probability plot do not appear to closely follow a straight line but show what appears to be a “wave”. We can compare these characteristics to the sample of 400 normally distributed observations in Example 3.24 and see that they represent much stronger deviations from the normal model. NBA player heights do not appear to come from a normal distribution.

- **Example 3.26** Can we approximate poker winnings by a normal distribution? We consider the poker winnings of an individual over 50 days. A histogram and normal probability plot of these data are shown in Figure 3.13.

The data are very strongly right skewed in the histogram, which corresponds to the very strong deviations on the upper right component of the normal probability plot. If we compare these results to the sample of 40 normal observations in Example 3.24, it is apparent that these data show very strong deviations from the normal model.

- **Guided Practice 3.27** Determine which data sets represented in Figure 3.14 plausibly come from a nearly normal distribution. Are you confident in all of your conclusions? There are 100 (top left), 50 (top right), 500 (bottom left), and 15 points (bottom right) in the four plots.<sup>21</sup>

<sup>20</sup>These data were collected from [www.nba.com](http://www.nba.com).

<sup>21</sup>Answers may vary a little. The top-left plot shows some deviations in the smallest values in the data set; specifically, the left tail of the data set has some outliers we should be wary of. The top-right and bottom-left plots do not show any obvious or extreme deviations from the lines for their respective sample sizes, so a normal model would be reasonable for these data sets. The bottom-right plot has a consistent curvature that suggests it is not from the normal distribution. If we examine just the vertical coordinates of these observations, we see that there is a lot of data between -20 and 0, and then about five observations scattered between 0 and 70. This describes a distribution that has a strong right skew.

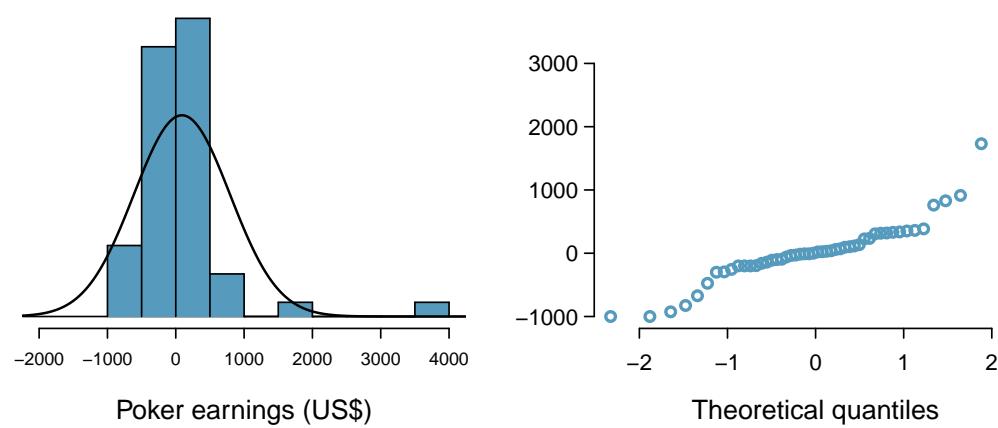


Figure 3.13: A histogram of poker data with the best fitting normal plot and a normal probability plot.

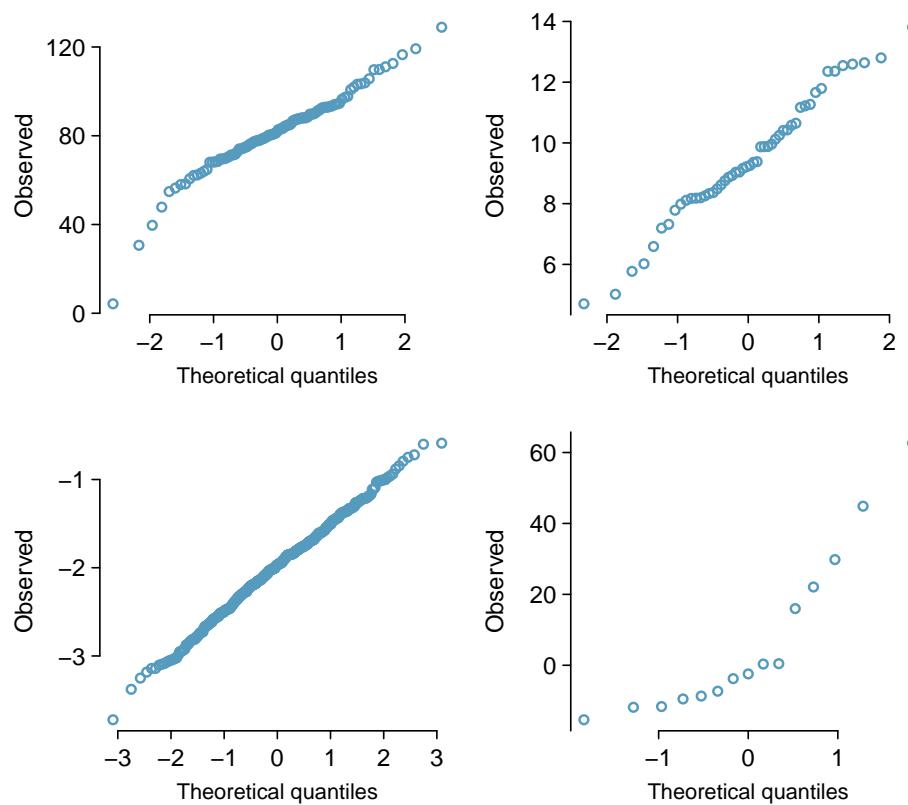


Figure 3.14: Four normal probability plots for Guided Practice 3.27.

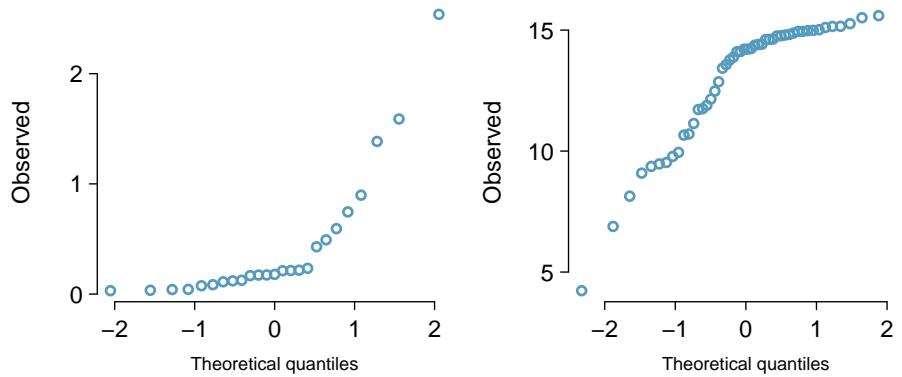


Figure 3.15: Normal probability plots for Guided Practice 3.28.

- Ⓐ **Guided Practice 3.28** Figure 3.15 shows normal probability plots for two distributions that are skewed. One distribution is skewed to the low end (left skewed) and the other to the high end (right skewed). Which is which?<sup>22</sup>

<sup>22</sup>Examine where the points fall along the vertical axis. In the first plot, most points are near the low end with fewer observations scattered along the high end; this describes a distribution that is skewed to the high end. The second plot shows the opposite features, and this distribution is skewed to the low end.

## 3.3 Exercises

### 3.3.1 Normal distribution

**3.1 Area under the curve, Part I.** What percent of a standard normal distribution  $N(\mu = 0, \sigma = 1)$  is found in each region? Be sure to draw a graph.

- (a)  $Z < -1.35$       (b)  $Z > 1.48$       (c)  $-0.4 < Z < 1.5$       (d)  $|Z| > 2$

**3.2 Area under the curve, Part II.** What percent of a standard normal distribution  $N(\mu = 0, \sigma = 1)$  is found in each region? Be sure to draw a graph.

- (a)  $Z > -1.13$       (b)  $Z < 0.18$       (c)  $Z > 8$       (d)  $|Z| < 0.5$

**3.3 GRE scores, Part I.** Sophia who took the Graduate Record Examination (GRE) scored 160 on the Verbal Reasoning section and 157 on the Quantitative Reasoning section. The mean score for Verbal Reasoning section for all test takers was 151 with a standard deviation of 7, and the mean score for the Quantitative Reasoning was 153 with a standard deviation of 7.67. Suppose that both distributions are nearly normal.

- (a) Write down the short-hand for these two normal distributions.
- (b) What is Sophia's Z-score on the Verbal Reasoning section? On the Quantitative Reasoning section? Draw a standard normal distribution curve and mark these two Z-scores.
- (c) What do these Z-scores tell you?
- (d) Relative to others, which section did she do better on?
- (e) Find her percentile scores for the two exams.
- (f) What percent of the test takers did better than her on the Verbal Reasoning section? On the Quantitative Reasoning section?
- (g) Explain why simply comparing raw scores from the two sections could lead to an incorrect conclusion as to which section a student did better on.
- (h) If the distributions of the scores on these exams are not nearly normal, would your answers to parts (b) - (f) change? Explain your reasoning.

**3.4 Triathlon times, Part I.** In triathlons, it is common for racers to be placed into age and gender groups. Friends Leo and Mary both completed the Hermosa Beach Triathlon, where Leo competed in the *Men, Ages 30 - 34* group while Mary competed in the *Women, Ages 25 - 29* group. Leo completed the race in 1:22:28 (4948 seconds), while Mary completed the race in 1:31:53 (5513 seconds). Obviously Leo finished faster, but they are curious about how they did within their respective groups. Can you help them? Here is some information on the performance of their groups:

- The finishing times of the *Men, Ages 30 - 34* group has a mean of 4313 seconds with a standard deviation of 583 seconds.
- The finishing times of the *Women, Ages 25 - 29* group has a mean of 5261 seconds with a standard deviation of 807 seconds.
- The distributions of finishing times for both groups are approximately Normal.

Remember: a better performance corresponds to a faster finish.

- (a) Write down the short-hand for these two normal distributions.
- (b) What are the Z-scores for Leo's and Mary's finishing times? What do these Z-scores tell you?
- (c) Did Leo or Mary rank better in their respective groups? Explain your reasoning.
- (d) What percent of the triathletes did Leo finish faster than in his group?
- (e) What percent of the triathletes did Mary finish faster than in her group?
- (f) If the distributions of finishing times are not nearly normal, would your answers to parts (b) - (e) change? Explain your reasoning.

**3.5 GRE scores, Part II.** In Exercise 3.3 we saw two distributions for GRE scores:  $N(\mu = 151, \sigma = 7)$  for the verbal part of the exam and  $N(\mu = 153, \sigma = 7.67)$  for the quantitative part. Use this information to compute each of the following:

- The score of a student who scored in the 80<sup>th</sup> percentile on the Quantitative Reasoning section.
- The score of a student who scored worse than 70% of the test takers in the Verbal Reasoning section.

**3.6 Triathlon times, Part II.** In Exercise 3.4 we saw two distributions for triathlon times:  $N(\mu = 4313, \sigma = 583)$  for *Men, Ages 30 - 34* and  $N(\mu = 5261, \sigma = 807)$  for the *Women, Ages 25 - 29* group. Times are listed in seconds. Use this information to compute each of the following:

- The cutoff time for the fastest 5% of athletes in the men's group, i.e. those who took the shortest 5% of time to finish.
- The cutoff time for the slowest 10% of athletes in the women's group.



**3.7 LA weather, Part I.** The average daily high temperature in June in LA is 77°F with a standard deviation of 5°F. Suppose that the temperatures in June closely follow a normal distribution.

- What is the probability of observing an 83°F temperature or higher in LA during a randomly chosen day in June?
- How cool are the coldest 10% of the days (days with lowest average high temperature) during June in LA?

**3.8 CAPM.** The Capital Asset Pricing Model (CAPM) is a financial model that assumes returns on a portfolio are normally distributed. Suppose a portfolio has an average annual return of 14.7% (i.e. an average gain of 14.7%) with a standard deviation of 33%. A return of 0% means the value of the portfolio doesn't change, a negative return means that the portfolio loses money, and a positive return means that the portfolio gains money.

- What percent of years does this portfolio lose money, i.e. have a return less than 0%?
- What is the cutoff for the highest 15% of annual returns with this portfolio?

**3.9 LA weather, Part II.** Exercise 3.7 states that average daily high temperature in June in LA is 77°F with a standard deviation of 5°F, and it can be assumed that they to follow a normal distribution. We use the following equation to convert °F (Fahrenheit) to °C (Celsius):

$$C = (F - 32) \times \frac{5}{9}.$$

- Write the probability model for the distribution of temperature in °C in June in LA.
- What is the probability of observing a 28°C (which roughly corresponds to 83°F) temperature or higher in June in LA? Calculate using the °C model from part (a).
- Did you get the same answer or different answers in part (b) of this question and part (a) of Exercise 3.7? Are you surprised? Explain.
- Estimate the IQR of the temperatures (in °C) in June in LA.

**3.10 Heights of 10 year olds.** Heights of 10 year olds, regardless of gender, closely follow a normal distribution with mean 55 inches and standard deviation 6 inches.

- What is the probability that a randomly chosen 10 year old is shorter than 48 inches?
- What is the probability that a randomly chosen 10 year old is between 60 and 65 inches?
- If the tallest 10% of the class is considered "very tall", what is the height cutoff for "very tall"?
- The height requirement for *Batman the Ride* at Six Flags Magic Mountain is 54 inches. What percent of 10 year olds cannot go on this ride?

**3.11 Auto insurance premiums.** Suppose a newspaper article states that the distribution of auto insurance premiums for residents of California is approximately normal with a mean of \$1,650. The article also states that 25% of California residents pay more than \$1,800.

- (a) What is the Z-score that corresponds to the top 25% (or the 75<sup>th</sup> percentile) of the standard normal distribution?
- (b) What is the mean insurance cost? What is the cutoff for the 75th percentile?
- (c) Identify the standard deviation of insurance premiums in California.

**3.12 Speeding on the I-5, Part I.** The distribution of passenger vehicle speeds traveling on the Interstate 5 Freeway (I-5) in California is nearly normal with a mean of 72.6 miles/hour and a standard deviation of 4.78 miles/hour.<sup>23</sup>

- (a) What percent of passenger vehicles travel slower than 80 miles/hour?
- (b) What percent of passenger vehicles travel between 60 and 80 miles/hour?
- (c) How fast do the fastest 5% of passenger vehicles travel?
- (d) The speed limit on this stretch of the I-5 is 70 miles/hour. Approximate what percentage of the passenger vehicles travel above the speed limit on this stretch of the I-5.

**3.13 Overweight baggage.** Suppose weights of the checked baggage of airline passengers follow a nearly normal distribution with mean 45 pounds and standard deviation 3.2 pounds. Most airlines charge a fee for baggage that weigh in excess of 50 pounds. Determine what percent of airline passengers incur this fee.

**3.14 Find the SD.** Find the standard deviation of the distribution in the following situations.

- (a) MENSA is an organization whose members have IQs in the top 2% of the population. IQs are normally distributed with mean 100, and the minimum IQ score required for admission to MENSA is 132.
- (b) Cholesterol levels for women aged 20 to 34 follow an approximately normal distribution with mean 185 milligrams per deciliter (mg/dl). Women with cholesterol levels above 220 mg/dl are considered to have high cholesterol and about 18.5% of women fall into this category.

**3.15 Buying books on Ebay.** The textbook you need to buy for your chemistry class is expensive at the college bookstore, so you consider buying it on Ebay instead. A look at past auctions suggest that the prices of that chemistry textbook have an approximately normal distribution with mean \$89 and standard deviation \$15.

- (a) What is the probability that a randomly selected auction for this book closes at more than \$100?
- (b) Ebay allows you to set your maximum bid price so that if someone outbids you on an auction you can automatically outbid them, up to the maximum bid price you set. If you are only bidding on one auction, what are the advantages and disadvantages of setting a bid price too high or too low? What if you are bidding on multiple auctions?
- (c) If you watched 10 auctions, roughly what percentile might you use for a maximum bid cutoff to be somewhat sure that you will win one of these ten auctions? Is it possible to find a cutoff point that will ensure that you win an auction?
- (d) If you are willing to track up to ten auctions closely, about what price might you use as your maximum bid price if you want to be somewhat sure that you will buy one of these ten books?

**3.16 SAT scores.** SAT scores (out of 2400) are distributed normally with a mean of 1500 and a standard deviation of 300. Suppose a school council awards a certificate of excellence to all students who score at least 1900 on the SAT, and suppose we pick one of the recognized students at random. What is the probability this student's score will be at least 2100? (The material covered in Section ?? would be useful for this question.)

---

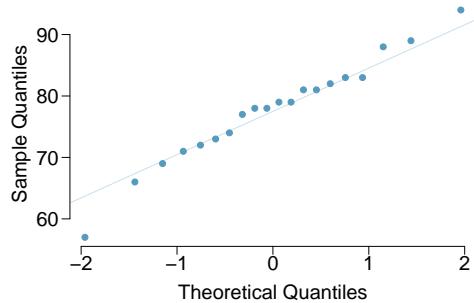
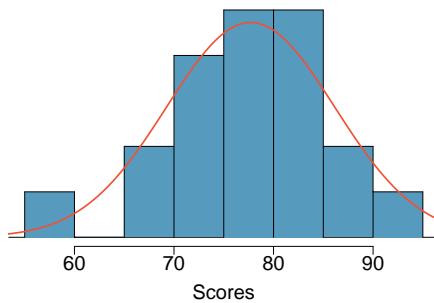
<sup>23</sup>S. Johnson and D. Murray. "Empirical Analysis of Truck and Automobile Speeds on Rural Interstates: Impact of Posted Speed Limits". In: *Transportation Research Board 89th Annual Meeting*. 2010.

### 3.3.2 Evaluating the normal approximation

**3.17 Scores on stats final.** Below are final exam scores of 20 Introductory Statistics students.

$$\begin{array}{cccccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 \\ 57, 66, 69, 71, 72, 73, 74, 77, 78, 78, 79, 79, 81, 81, 82, 83, 83, 88, 89, 94 \end{array}$$

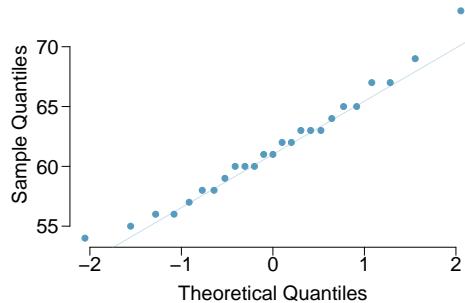
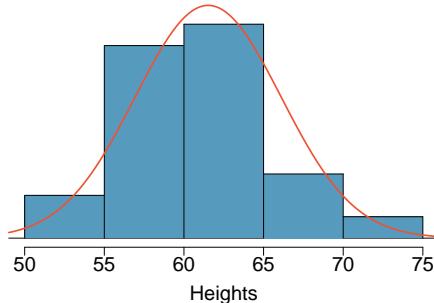
- (a) The mean score is 77.7 points, with a standard deviation of 8.44 points. Use this information to determine if the scores approximately follow the 68-95-99.7% Rule.
- (b) Do these data appear to follow a normal distribution? Explain your reasoning using the graphs provided below.



**3.18 Heights of female college students.** Below are heights of 25 female college students.

$$\begin{array}{cccccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 & 21 & 22 & 23 & 24 & 25 \\ 54, 55, 56, 56, 57, 57, 58, 58, 59, 60, 60, 60, 61, 61, 62, 62, 63, 63, 63, 64, 65, 65, 67, 67, 69, 73 \end{array}$$

- (a) The mean height is 61.52 inches with a standard deviation of 4.58 inches. Use this information to determine if the heights approximately follow the 68-95-99.7% Rule.
- (b) Do these data appear to follow a normal distribution? Explain your reasoning using the graphs provided below.



# Chapter 4

## Foundations for inference

Statistical inference is concerned primarily with understanding the quality of parameter estimates. For example, a classic inferential question is, “How sure are we that the estimated mean,  $\bar{x}$ , is near the true population mean,  $\mu$ ? ” While the equations and details change depending on the setting, the foundations for inference are the same throughout all of statistics. We introduce these common themes in Sections 4.1-4.4 by discussing inference about the population mean,  $\mu$ , and set the stage for other parameters and scenarios in Section ???. Understanding this chapter will make the rest of this book, and indeed the rest of statistics, seem much more familiar.

Throughout the next few sections we consider a data set called `yrbss`, which represents all 13,583 high school students in the Youth Risk Behavior Surveillance System (YRBSS) from 2013.<sup>1</sup> Part of this data set is shown in Table 4.1, and the variables are described in Table 4.2.

ID	age	gender	grade	height	weight	helmet	active	lifting
1	14	female	9			never	4	0
2	14	female	9			never	2	0
3	15	female	9	1.73	84.37	never	7	0
:	:	:	:	:	:	:	:	:
13582	17	female	12	1.60	77.11	sometimes	5	
13583	17	female	12	1.57	52.16	did not ride	5	

Table 4.1: Five cases from the `yrbss` data set. Some observations are blank since there are missing data. For example, the height and weight of students 1 and 2 are missing.

We’re going to consider the population of high school students who participated in the 2013 YRBSS. We took a simple random sample of this population, which is represented in Table 4.3.<sup>2</sup> We will use this sample, which we refer to as the `yrbss_samp` data set, to draw conclusions about the population of YRBSS participants. This is the practice of statistical inference in the broadest sense. Four histograms summarizing the `height`, `weight`, `active`, and `lifting` variables from `yrbss_samp` data set are shown in Figure 4.4.

<sup>1</sup>[www.cdc.gov/healthyyouth/data/yrbs/data.htm](http://www.cdc.gov/healthyyouth/data/yrbs/data.htm)

<sup>2</sup>About 10% of high schoolers for each variable chose not to answer the question, we used multiple regression (see Chapter ???) to predict what those responses would have been. For simplicity, we will assume that these predicted values are exactly the truth.

Name	Description
age	Age of the student.
gender	Sex of the student.
grade	Grade in high school
height	Height, in meters. There are 3.28 feet in a meter.
weight	Weight, in kilograms (2.2 pounds per kilogram).
helmet	Frequency that the student wore a helmet while biking in the last 12 months.
active	Number of days physically active for 60+ minutes in the last 7 days.
lifting	Number of days of strength training (e.g. lifting weights) in the last 7 days.

Table 4.2: Variables and their descriptions for the `yrbss` data set.

ID	age	gender	grade	height	weight	helmet	active	lifting
5653	16	female	11	1.50	52.62	never	0	0
9437	17	male	11	1.78	74.84	rarely	7	5
2021	17	male	11	1.75	106.60	never	7	0
:	:	:	:	:	:	:	:	:
2325	14	male	9	1.70	55.79	never	1	0

Table 4.3: Four observations for the `yrbss_samp` data set, which represents a simple random sample of 100 high schoolers from the 2013 YRBSS.

## 4.1 Variability in estimates

We would like to estimate four features of the high schoolers in YRBSS using the sample.

- (1) What is the average height of the YRBSS high schoolers?
- (2) What is the average weight of the YRBSS high schoolers?
- (3) On average, how many days per week are YRBSS high schoolers physically active?
- (4) On average, how many days per week do YRBSS high schoolers do weight training?

While we focus on the mean in this chapter, questions regarding variation are often just as important in practice. For instance, if students are either very active or almost entirely inactive (the distribution is bimodal), we might try different strategies to promote a healthy lifestyle among students than if all high schoolers were already somewhat active.

### 4.1.1 Point estimates

We want to estimate the **population mean** based on the sample. The most intuitive way to go about doing this is to simply take the **sample mean**. That is, to estimate the average height of all YRBSS students, take the average height for the sample:

$$\bar{x}_{\text{height}} = \frac{1.50 + 1.78 + \dots + 1.70}{100} = 1.697$$

The sample mean  $\bar{x} = 1.697$  meters (5 feet, 6.8 inches) is called a **point estimate** of the population mean: if we can only choose one value to estimate the population mean, this is our best guess. Suppose we take a new sample of 100 people and recompute the mean; we will probably not get the exact same answer that we got using the `yrbss_samp` data set. Estimates generally vary from one sample to another, and this **sampling variation** suggests our estimate may be close, but it will not be exactly equal to the parameter.

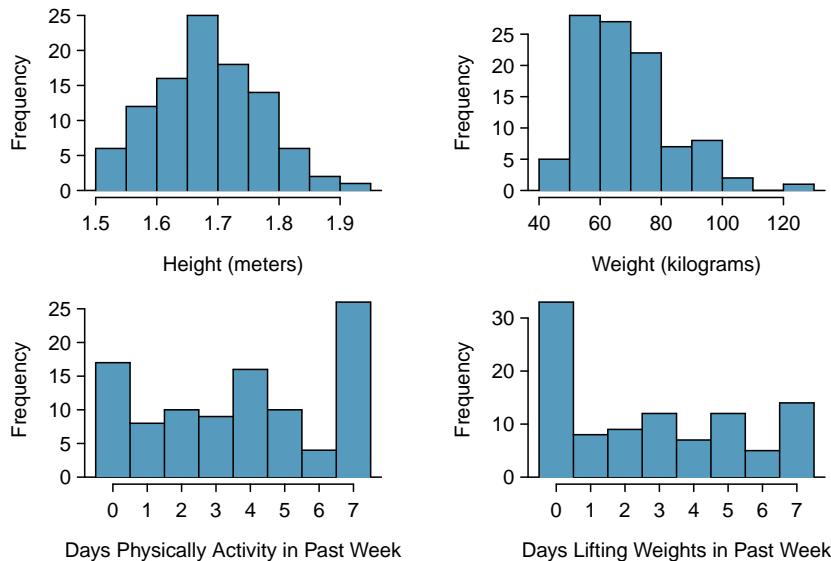


Figure 4.4: Histograms of `height`, `weight`, `active`, and `lifting` for the sample YRBSS data. The `height` distribution is approximately symmetric, `weight` is moderately skewed to the right, `activity` is bimodal or multimodal (with unclear skew), and `lifting` is strongly right skewed.

We can also estimate the average weight of YRBSS respondents by examining the sample mean of `weight` (in kg), and average number of days physically active in a week:

$$\bar{x}_{\text{weight}} = \frac{52.6 + 74.8 + \dots + 55.8}{100} = 68.89 \quad \bar{x}_{\text{active}} = \frac{0 + 7 + \dots + 1}{100} = 3.75$$

The average weight is 68.89 kilograms, which is about 151.6 pounds.

What about generating point estimates of other **population parameters**, such as the population median or population standard deviation? Once again we might estimate parameters based on sample statistics, as shown in Table 4.5. For example, the population standard deviation of `active` using the sample standard deviation, 2.56 days.

active	estimate	parameter
mean	3.75	3.90
median	4.00	4.00
st. dev.	2.556	2.564

Table 4.5: Point estimates and parameter values for the `active` variable. The parameters were obtained by computing the mean, median, and SD for all YRBSS respondents.

- Ⓐ **Guided Practice 4.1** Suppose we want to estimate the difference in days active for men and women. If  $\bar{x}_{\text{men}} = 4.3$  and  $\bar{x}_{\text{women}} = 3.2$ , then what would be a good point estimate for the population difference?<sup>3</sup>

<sup>3</sup>We could take the difference of the two sample means:  $4.3 - 3.2 = 1.1$ . Men are physically active about 1.1 days per week more than women on average in YRBSS.

Ⓐ **Guided Practice 4.2** If you had to provide a point estimate of the population IQR for the heights of participants, how might you make such an estimate using a sample?<sup>4</sup>

### 4.1.2 Point estimates are not exact

Estimates are usually not exactly equal to the truth, but they get better as more data become available. We can see this by plotting a running mean from `yrbss_samp`. A **running mean** is a sequence of means, where each mean uses one more observation in its calculation than the mean directly before it in the sequence. For example, the second mean in the sequence is the average of the first two observations and the third in the sequence is the average of the first three. The running mean for the `active` variable in the `yrbss_samp` is shown in Figure 4.6, and it approaches the true population average, 3.90 days, as more data become available.

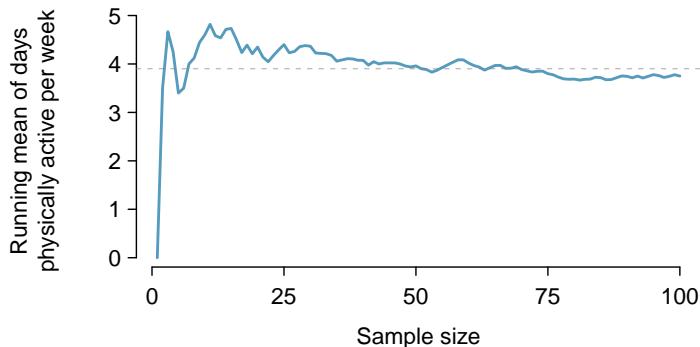


Figure 4.6: The mean computed after adding each individual to the sample. The mean tends to approach the true population average as more data become available.

Sample point estimates only approximate the population parameter, and they vary from one sample to another. If we took another simple random sample of the YRBSS students, we would find that the sample mean for the number of days active would be a little different. It will be useful to quantify how variable an estimate is from one sample to another. If this variability is small (i.e. the sample mean doesn't change much from one sample to another) then that estimate is probably very accurate. If it varies widely from one sample to another, then we should not expect our estimate to be very good.

### 4.1.3 Standard error of the mean

From the random sample represented in `yrbss_samp`, we guessed the average number of days a YRBSS student is physically active is 3.75 days. Suppose we take another random sample of 100 individuals and take its mean: 3.22 days. Suppose we took another (3.67 days) and another (4.10 days), and so on. If we do this many many times – which we can do only because we have all YRBSS students – we can build up a **sampling distribution** for the sample mean when the sample size is 100, shown in Figure 4.7.

---

<sup>4</sup>To obtain a point estimate of the height for the full set of YRBSS students, we could take the IQR of the sample.

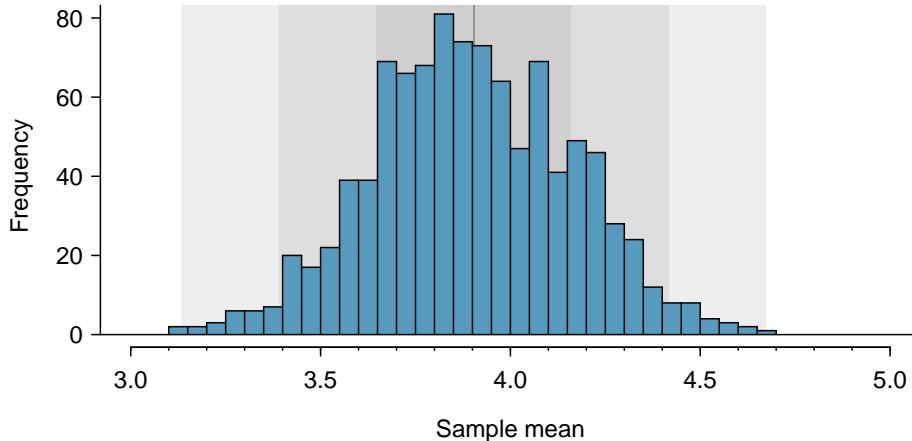


Figure 4.7: A histogram of 1000 sample means for number of days physically active per week, where the samples are of size  $n = 100$ .

### Sampling distribution

The sampling distribution represents the distribution of the point estimates based on samples of a fixed size from a certain population. It is useful to think of a particular point estimate as being drawn from such a distribution. Understanding the concept of a sampling distribution is central to understanding statistical inference.

The sampling distribution shown in Figure 4.7 is unimodal and approximately symmetric. It is also centered exactly at the true population mean:  $\mu = 3.90$ . Intuitively, this makes sense. The sample means should tend to “fall around” the population mean.

We can see that the sample mean has some variability around the population mean, which can be quantified using the standard deviation of this distribution of sample means:  $\sigma_{\bar{x}} = 0.26$ . The standard deviation of the sample mean tells us how far the typical estimate is away from the actual population mean, 3.90 days. It also describes the typical **error** of the point estimate, and for this reason we usually call this standard deviation the **standard error (SE)** of the estimate.

*SE*  
standard  
error

### Standard error of an estimate

The standard deviation associated with an estimate is called the *standard error*. It describes the typical error or uncertainty associated with the estimate.

When considering the case of the point estimate  $\bar{x}$ , there is one problem: there is no obvious way to estimate its standard error from a single sample. However, statistical theory provides a helpful tool to address this issue.

- Ⓐ **Guided Practice 4.3** (a) Would you rather use a small sample or a large sample when estimating a parameter? Why? (b) Using your reasoning from (a), would you expect a point estimate based on a small sample to have smaller or larger standard error than a point estimate based on a larger sample?<sup>5</sup>

In the sample of 100 students, the standard error of the sample mean is equal to the population standard deviation divided by the square root of the sample size:

$$SE_{\bar{x}} = \sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}} = \frac{2.6}{\sqrt{100}} = 0.26$$

where  $\sigma_x$  is the standard deviation of the individual observations. This is no coincidence. We can show mathematically that this equation is correct when the observations are independent using the probability tools of Section ??.

#### Computing SE for the sample mean

Given  $n$  independent observations from a population with standard deviation  $\sigma$ , the standard error of the sample mean is equal to

$$SE = \frac{\sigma}{\sqrt{n}} \tag{4.4}$$

A reliable method to ensure sample observations are independent is to conduct a simple random sample consisting of less than 10% of the population.

There is one subtle issue in Equation (4.4): the population standard deviation is typically unknown. You might have already guessed how to resolve this problem: we can use the point estimate of the standard deviation from the sample. This estimate tends to be sufficiently good when the sample size is at least 30 and the population distribution is not strongly skewed. Thus, we often just use the sample standard deviation  $s$  instead of  $\sigma$ . When the sample size is smaller than 30, we will need to use a method to account for extra uncertainty in the standard error. If the skew condition is not met, a larger sample is needed to compensate for the extra skew. These topics are further discussed in Section 4.4.

- Ⓐ **Guided Practice 4.5** In the sample of 100 students, the standard deviation of student heights is  $s_{height} = 0.088$  meters. In this case, we can confirm that the observations are independent by checking that the data come from a simple random sample consisting of less than 10% of the population. (a) What is the standard error of the sample mean,  $\bar{x}_{height} = 1.70$  meters? (b) Would you be surprised if someone told you the average height of all YRBSS respondents was actually 1.69 meters?<sup>6</sup>

---

<sup>5</sup>(a) Consider two random samples: one of size 10 and one of size 1000. Individual observations in the small sample are highly influential on the estimate while in larger samples these individual observations would more often average each other out. The larger sample would tend to provide a more accurate estimate. (b) If we think an estimate is better, we probably mean it typically has less error. Based on (a), our intuition suggests that a larger sample size corresponds to a smaller standard error.

<sup>6</sup>(a) Use Equation (4.4) with the sample standard deviation to compute the standard error:  $SE_{\bar{y}} = 0.088/\sqrt{100} = 0.0088$  meters. (b) It would not be surprising. Our sample is about 1 standard error from 1.69m. In other words, 1.69m does not seem to be implausible given that our sample was relatively close to it. (We use the standard error to identify what is close.)

Ⓐ **Guided Practice 4.6** (a) Would you be more trusting of a sample that has 100 observations or 400 observations? (b) We want to show mathematically that our estimate tends to be better when the sample size is larger. If the standard deviation of the individual observations is 10, what is our estimate of the standard error when the sample size is 100? What about when it is 400? (c) Explain how your answer to part (b) mathematically justifies your intuition in part (a).<sup>7</sup>

#### 4.1.4 Basic properties of point estimates

We achieved three goals in this section. First, we determined that point estimates from a sample may be used to estimate population parameters. We also determined that these point estimates are not exact: they vary from one sample to another. Lastly, we quantified the uncertainty of the sample mean using what we call the standard error, mathematically represented in Equation (4.4). While we could also quantify the standard error for other estimates – such as the median, standard deviation, or any other number of statistics – we will postpone these extensions until later chapters or courses.

## 4.2 Confidence intervals

A point estimate provides a single plausible value for a parameter. However, a point estimate is rarely perfect; usually there is some error in the estimate. Instead of supplying just a point estimate of a parameter, a next logical step would be to provide a plausible *range of values* for the parameter.

### 4.2.1 Capturing the population parameter

A plausible range of values for the population parameter is called a **confidence interval**.

Using only a point estimate is like fishing in a murky lake with a spear, and using a confidence interval is like fishing with a net. We can throw a spear where we saw a fish, but we will probably miss. On the other hand, if we toss a net in that area, we have a good chance of catching the fish.

If we report a point estimate, we probably will not hit the exact population parameter. On the other hand, if we report a range of plausible values – a confidence interval – we have a good shot at capturing the parameter.

Ⓐ **Guided Practice 4.7** If we want to be very certain we capture the population parameter, should we use a wider interval or a smaller interval?<sup>8</sup>

---

<sup>7</sup>(a) Extra observations are usually helpful in understanding the population, so a point estimate with 400 observations seems more trustworthy. (b) The standard error when the sample size is 100 is given by  $SE_{100} = 10/\sqrt{100} = 1$ . For 400:  $SE_{400} = 10/\sqrt{400} = 0.5$ . The larger sample has a smaller standard error. (c) The standard error of the sample with 400 observations is lower than that of the sample with 100 observations. The standard error describes the typical error, and since it is lower for the larger sample, this mathematically shows the estimate from the larger sample tends to be better – though it does not guarantee that every large sample will provide a better estimate than a particular small sample.

<sup>8</sup>If we want to be more certain we will capture the fish, we might use a wider net. Likewise, we use a wider confidence interval if we want to be more certain that we capture the parameter.

### 4.2.2 An approximate 95% confidence interval

Our point estimate is the most plausible value of the parameter, so it makes sense to build the confidence interval around the point estimate. The standard error, which is a measure of the uncertainty associated with the point estimate, provides a guide for how large we should make the confidence interval.

The standard error represents the standard deviation associated with the estimate, and roughly 95% of the time the estimate will be within 2 standard errors of the parameter. If the interval spreads out 2 standard errors from the point estimate, we can be roughly 95% **confident** that we have captured the true parameter:

$$\text{point estimate} \pm 2 \times SE \quad (4.8)$$

But what does “95% confident” mean? Suppose we took many samples and built a confidence interval from each sample using Equation (4.8). Then about 95% of those intervals would contain the actual mean,  $\mu$ . Figure 4.8 shows this process with 25 samples, where 24 of the resulting confidence intervals contain the average number of days per week that YRBSS students are physically active,  $\mu = 3.90$  days, and one interval does not.

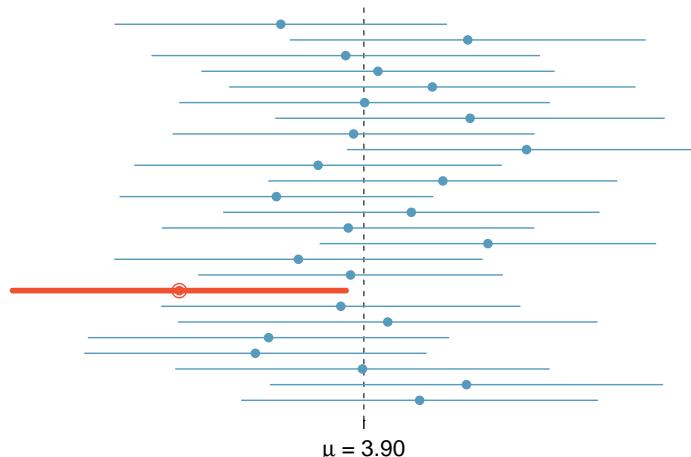


Figure 4.8: Twenty-five samples of size  $n = 100$  were taken from `yrbss`. For each sample, a confidence interval was created to try to capture the average number of days per week that students are physically active. Only 1 of these 25 intervals did not capture the true mean,  $\mu = 3.90$  days.

- **Guided Practice 4.9** In Figure 4.8, one interval does not contain 3.90 days. Does this imply that the mean cannot be 3.90?<sup>9</sup>

The rule where about 95% of observations are within 2 standard deviations of the mean is only approximately true. However, it holds very well for the normal distribution. As we will soon see, the mean tends to be normally distributed when the sample size is sufficiently large.

---

<sup>9</sup>Just as some observations occur more than 2 standard deviations from the mean, some point estimates will be more than 2 standard errors from the parameter. A confidence interval only provides a plausible range of values for a parameter. While we might say other values are implausible based on the data, this does not mean they are impossible.

● **Example 4.10** The sample mean of days active per week from `yrbss_samp` is 3.75 days. The standard error, as estimated using the sample standard deviation, is  $SE = \frac{2.6}{\sqrt{100}} = 0.26$  days. (The population SD is unknown in most applications, so we use the sample SD here.) Calculate an approximate 95% confidence interval for the average days active per week for all YRBSS students.

We apply Equation (4.8):

$$3.75 \pm 2 \times 0.26 \rightarrow (3.23, 4.27)$$

Based on these data, we are about 95% confident that the average days active per week for all YRBSS students was larger than 3.23 but less than 4.27 days. Our interval extends out 2 standard errors from the point estimate,  $\bar{x}_{active}$ .

○ **Guided Practice 4.11** The sample data suggest the average YRBSS student height is  $\bar{x}_{height} = 1.697$  meters with a standard error of 0.0088 meters (estimated using the sample standard deviation, 0.088 meters). What is an approximate 95% confidence interval for the average height of all of the YRBSS students?<sup>10</sup>

### 4.2.3 The sampling distribution for the mean

In Section 4.1.3, we introduced a sampling distribution for  $\bar{x}$ , the average days physically active per week for samples of size 100. We examined this distribution earlier in Figure 4.7. Now we'll take 100,000 samples, calculate the mean of each, and plot them in a histogram to get an especially accurate depiction of the sampling distribution. This histogram is shown in the left panel of Figure 4.9.

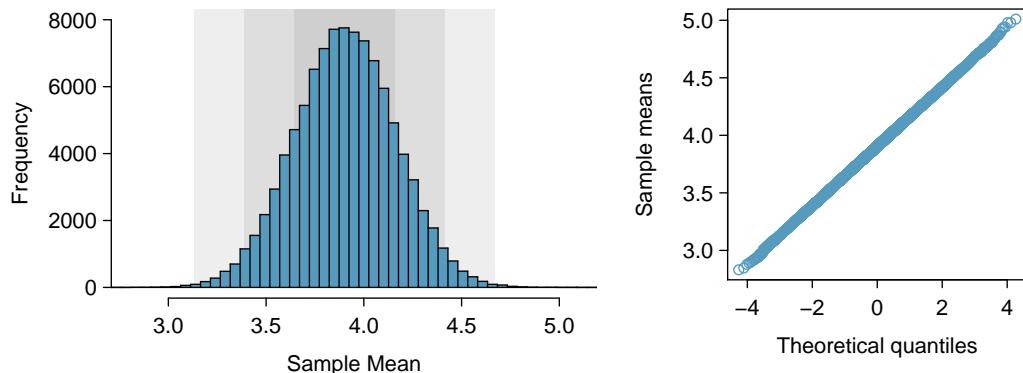


Figure 4.9: The left panel shows a histogram of the sample means for 100,000 different random samples. The right panel shows a normal probability plot of those sample means.

Does this distribution look familiar? Hopefully so! The distribution of sample means closely resembles the normal distribution (see Section 3.1). A normal probability plot of these sample means is shown in the right panel of Figure 4.9. Because all of the points

<sup>10</sup>Apply Equation (4.8):  $1.697 \pm 2 \times 0.0088 \rightarrow (1.6794, 1.7146)$ . We interpret this interval as follows: We are about 95% confident the average height of all YRBSS students was between 1.6794 and 1.7146 meters (5.51 to 5.62 feet).

closely fall around a straight line, we can conclude the distribution of sample means is nearly normal. This result can be explained by the Central Limit Theorem.

#### Central Limit Theorem, informal description

If a sample consists of at least 30 independent observations and the data are not strongly skewed, then the distribution of the sample mean is well approximated by a normal model.

We will apply this informal version of the Central Limit Theorem for now, and discuss its details further in Section 4.4.

The choice of using 2 standard errors in Equation (4.8) was based on our general guideline that roughly 95% of the time, observations are within two standard deviations of the mean. Under the normal model, we can make this more accurate by using 1.96 in place of 2.

$$\text{point estimate} \pm 1.96 \times SE \quad (4.12)$$

If a point estimate, such as  $\bar{x}$ , is associated with a normal model and standard error  $SE$ , then we use this more precise 95% confidence interval.

#### 4.2.4 Changing the confidence level

Suppose we want to consider confidence intervals where the confidence level is somewhat higher than 95%; perhaps we would like a confidence level of 99%. Think back to the analogy about trying to catch a fish: if we want to be more sure that we will catch the fish, we should use a wider net. To create a 99% confidence level, we must also widen our 95% interval. On the other hand, if we want an interval with lower confidence, such as 90%, we could make our original 95% interval slightly slimmer.

The 95% confidence interval structure provides guidance in how to make intervals with new confidence levels. Below is a general 95% confidence interval for a point estimate that comes from a nearly normal distribution:

$$\text{point estimate} \pm 1.96 \times SE \quad (4.13)$$

There are three components to this interval: the point estimate, “1.96”, and the standard error. The choice of  $1.96 \times SE$  was based on capturing 95% of the data since the estimate is within 1.96 standard deviations of the parameter about 95% of the time. The choice of 1.96 corresponds to a 95% confidence level.

- **Guided Practice 4.14** If  $X$  is a normally distributed random variable, how often will  $X$  be within 2.58 standard deviations of the mean?<sup>11</sup>

To create a 99% confidence interval, change 1.96 in the 95% confidence interval formula to be 2.58. Guided Practice 4.14 highlights that 99% of the time a normal random variable will be within 2.58 standard deviations of the mean. This approach – using the Z-scores in the normal model to compute confidence levels – is appropriate when  $\bar{x}$  is associated with

<sup>11</sup>This is equivalent to asking how often the Z-score will be larger than -2.58 but less than 2.58. (For a picture, see Figure 4.10.) To determine this probability, look up -2.58 and 2.58 in the normal probability table (0.0049 and 0.9951). Thus, there is a  $0.9951 - 0.0049 \approx 0.99$  probability that the unobserved random variable  $X$  will be within 2.58 standard deviations of  $\mu$ .

a normal distribution with mean  $\mu$  and standard deviation  $SE_{\bar{x}}$ . Thus, the formula for a 99% confidence interval is

$$\bar{x} \pm 2.58 \times SE_{\bar{x}} \quad (4.15)$$

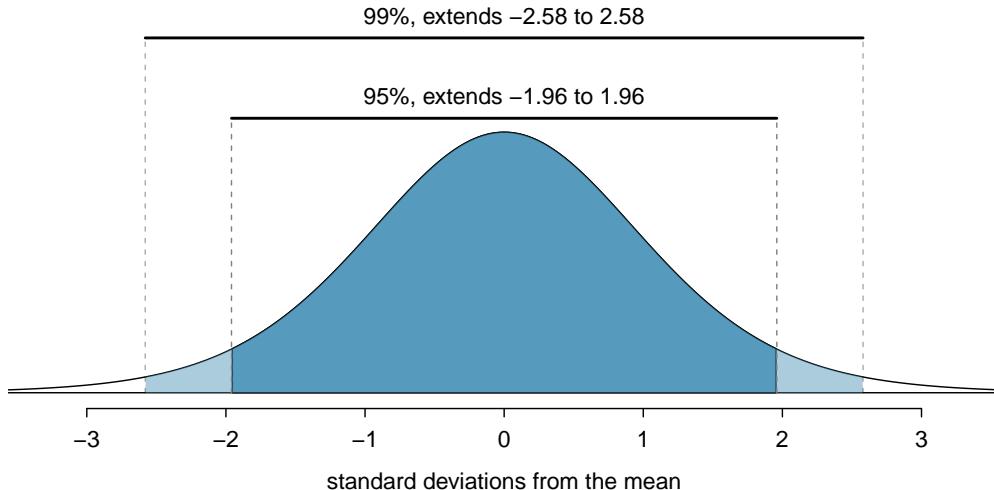


Figure 4.10: The area between  $-z^*$  and  $z^*$  increases as  $|z^*|$  becomes larger. If the confidence level is 99%, we choose  $z^*$  such that 99% of the normal curve is between  $-z^*$  and  $z^*$ , which corresponds to 0.5% in the lower tail and 0.5% in the upper tail:  $z^* = 2.58$ .

The normal approximation is crucial to the precision of these confidence intervals. Section 4.4 provides a more detailed discussion about when the normal model can safely be applied. When the normal model is not a good fit, we will use alternative distributions that better characterize the sampling distribution.

#### Conditions for $\bar{x}$ being nearly normal and $SE$ being accurate

Important conditions to help ensure the sampling distribution of  $\bar{x}$  is nearly normal and the estimate of  $SE$  sufficiently accurate:

- The sample observations are independent.
- The sample size is large:  $n \geq 30$  is a good rule of thumb.
- The population distribution is not strongly skewed. This condition can be difficult to evaluate, so just use your best judgement.

Additionally, the larger the sample size, the more lenient we can be with the sample's skew.

**How to verify sample observations are independent**

If the observations are from a simple random sample and consist of fewer than 10% of the population, then they are independent.

Subjects in an experiment are considered independent if they undergo random assignment to the treatment groups.

If a sample is from a seemingly random process, e.g. the lifetimes of wrenches used in a particular manufacturing process, checking independence is more difficult. In this case, use your best judgement.

**Checking for strong skew usually means checking for obvious outliers**

When there are prominent outliers present, the sample should contain at least 100 observations, and in some cases, much more.

This is a first course in statistics, so you won't have perfect judgement on assessing skew. That's okay. If you're in a bind, either consult a statistician or learn about the studentized bootstrap (bootstrap-t) method.

- **Guided Practice 4.16** Create a 99% confidence interval for the average days active per week of all YRBSS students using `yrbss_samp`. The point estimate is  $\bar{x}_{active} = 3.75$  and the standard error is  $SE_{\bar{x}} = 0.26$ .<sup>12</sup>

**Confidence interval for any confidence level**

If the point estimate follows the normal model with standard error  $SE$ , then a confidence interval for the population parameter is

$$\text{point estimate} \pm z^* SE$$

where  $z^*$  corresponds to the confidence level selected.

Figure 4.10 provides a picture of how to identify  $z^*$  based on a confidence level. We select  $z^*$  so that the area between  $-z^*$  and  $z^*$  in the normal model corresponds to the confidence level.

**Margin of error**

In a confidence interval,  $z^* \times SE$  is called the **margin of error**.

---

<sup>12</sup>The observations are independent (simple random sample, < 10% of the population), the sample size is at least 30 ( $n = 100$ ), and the distribution doesn't have a clear skew (Figure 4.4 on page 97); the normal approximation and estimate of SE should be reasonable. Apply the 99% confidence interval formula:  $\bar{x}_{active} \pm 2.58 \times SE_{\bar{x}} \rightarrow (3.08, 4.42)$ . We are 99% confident that the average days active per week of all YRBSS students is between 3.08 and 4.42 days.

Ⓐ **Guided Practice 4.17** Use the data in Guided Practice 4.16 to create a 90% confidence interval for the average days active per week of all YRBSS students.<sup>13</sup>

### 4.2.5 Interpreting confidence intervals

A careful eye might have observed the somewhat awkward language used to describe confidence intervals. Correct interpretation:

We are XX% confident that the population parameter is between...

*Incorrect* language might try to describe the confidence interval as capturing the population parameter with a certain probability. This is a common error: while it might be useful to think of it as a probability, the confidence level only quantifies how plausible it is that the parameter is in the interval.

Another important consideration of confidence intervals is that they *only try to capture the population parameter*. A confidence interval says nothing about the confidence of capturing individual observations, a proportion of the observations, or about capturing point estimates. Confidence intervals only attempt to capture population parameters.

## 4.3 Hypothesis testing

Are students lifting weights or performing other strength training exercises more or less often than they have in the past? We'll compare data from students from the 2011 YRBSS survey to our sample of 100 students from the 2013 YRBSS survey.

We'll also consider sleep behavior. A recent study found that college students average about 7 hours of sleep per night.<sup>14</sup> However, researchers at a rural college are interested in showing that their students sleep longer than seven hours on average. We investigate this topic in Section 4.3.4.

### 4.3.1 Hypothesis testing framework

Students from the 2011 YRBSS lifted weights (or performed other strength training exercises) 3.09 days per week on average. We want to determine if the `yrbss_samp` data set provides strong evidence that YRBSS students selected in 2013 are lifting more or less than the 2011 YRBSS students, versus the other possibility that there has been no change.<sup>15</sup> We simplify these three options into two competing **hypotheses**:

$H_0$ : The average days per week that YRBSS students lifted weights was the same for 2011 and 2013.

$H_A$ : The average days per week that YRBSS students lifted weights was *different* for 2013 than in 2011.

We call  $H_0$  the null hypothesis and  $H_A$  the alternative hypothesis.

---

<sup>13</sup>We first find  $z^*$  such that 90% of the distribution falls between  $-z^*$  and  $z^*$  in the standard normal model,  $N(\mu = 0, \sigma = 1)$ . We can look up  $-z^*$  in the normal probability table by looking for a lower tail of 5% (the other 5% is in the upper tail):  $z^* = 1.65$ . The 90% confidence interval can then be computed as  $\bar{x}_{active} \pm 1.65 \times SE_{\bar{x}} \rightarrow (3.32, 4.18)$ . (We had already verified conditions for normality and the standard error.) That is, we are 90% confident the average days active per week is between 3.32 and 4.18 days.

<sup>14</sup>Poll shows college students get least amount of sleep. [theloquitur.com/?p=1161](http://theloquitur.com/?p=1161)

<sup>15</sup>While we could answer this question by examining the entire YRBSS data set from 2013 (`yrbss`), we only consider the sample data (`yrbss_samp`), which is more realistic since we rarely have access to population data.

$H_0$   
null hypothesis

$H_A$   
alternative hypothesis

### Null and alternative hypotheses

The **null hypothesis** ( $H_0$ ) often represents either a skeptical perspective or a claim to be tested. The **alternative hypothesis** ( $H_A$ ) represents an alternative claim under consideration and is often represented by a range of possible parameter values.

The null hypothesis often represents a skeptical position or a perspective of no difference. The alternative hypothesis often represents a new perspective, such as the possibility that there has been a change.

### TIP: Hypothesis testing framework

The skeptic will not reject the null hypothesis ( $H_0$ ), unless the evidence in favor of the alternative hypothesis ( $H_A$ ) is so strong that she rejects  $H_0$  in favor of  $H_A$ .

The hypothesis testing framework is a very general tool, and we often use it without a second thought. If a person makes a somewhat unbelievable claim, we are initially skeptical. However, if there is sufficient evidence that supports the claim, we set aside our skepticism and reject the null hypothesis in favor of the alternative. The hallmarks of hypothesis testing are also found in the US court system.

- **Guided Practice 4.18** A US court considers two possible claims about a defendant: she is either innocent or guilty. If we set these claims up in a hypothesis framework, which would be the null hypothesis and which the alternative?<sup>16</sup>

Jurors examine the evidence to see whether it convincingly shows a defendant is guilty. Even if the jurors leave unconvinced of guilt beyond a reasonable doubt, this does not mean they believe the defendant is innocent. This is also the case with hypothesis testing: *even if we fail to reject the null hypothesis, we typically do not accept the null hypothesis as true.* Failing to find strong evidence for the alternative hypothesis is not equivalent to accepting the null hypothesis.

In the example with the YRBSS, the null hypothesis represents no difference in the average days per week of weight lifting in 2011 and 2013. The alternative hypothesis represents something new or more interesting: there was a difference, either an increase or a decrease. These hypotheses can be described in mathematical notation using  $\mu_{13}$  as the average days of weight lifting for 2013:

$$H_0: \mu_{13} = 3.09$$

$$H_A: \mu_{13} \neq 3.09$$

where 3.09 is the average number of days per week that students from the 2011 YRBSS lifted weights. Using the mathematical notation, the hypotheses can more easily be evaluated using statistical tools. We call 3.09 the **null value** since it represents the value of the parameter if the null hypothesis is true.

<sup>16</sup>The jury considers whether the evidence is so convincing (strong) that there is no reasonable doubt regarding the person's guilt; in such a case, the jury rejects innocence (the null hypothesis) and concludes the defendant is guilty (alternative hypothesis).

### 4.3.2 Testing hypotheses using confidence intervals

We will use the `yrbss_samp` data set to evaluate the hypothesis test, and we start by comparing the 2013 point estimate of the number of days per week that students lifted weights:  $\bar{x}_{13} = 2.78$  days. This estimate suggests that students from the 2013 YRBSS were lifting weights less than students in the 2011 YRBSS. However, to evaluate whether this provides strong evidence that there has been a change, we must consider the uncertainty associated with  $\bar{x}_{13}$ .

We learned in Section 4.1 that there is fluctuation from one sample to another, and it is unlikely that the sample mean will be exactly equal to the parameter; we should not expect  $\bar{x}_{13}$  to exactly equal  $\mu_{13}$ . Given that  $\bar{x}_{13} = 2.78$ , it might still be possible that the average of all students from the 2013 YRBSS survey is the same as the average from the 2011 YRBSS survey. The difference between  $\bar{x}_{13}$  and 3.09 could be due to *sampling variation*, i.e. the variability associated with the point estimate when we take a random sample.

In Section 4.2, confidence intervals were introduced as a way to find a range of plausible values for the population mean.

- **Example 4.19** In the sample of 100 students from the 2013 YRBSS survey, the average number of days per week that students lifted weights was 2.78 days with a standard deviation of 2.56 days (coincidentally the same as days active). Compute a 95% confidence interval for the average for all students from the 2013 YRBSS survey. You can assume the conditions for the normal model are met.

The general formula for the confidence interval based on the normal distribution is

$$\bar{x} \pm z^* SE_{\bar{x}}$$

We are given  $\bar{x}_{13} = 2.78$ , we use  $z^* = 1.96$  for a 95% confidence level, and we can compute the standard error using the standard deviation divided by the square root of the sample size:

$$SE_{\bar{x}} = \frac{s_{13}}{\sqrt{n}} = \frac{2.56}{\sqrt{100}} = 0.256$$

Entering the sample mean,  $z^*$ , and the standard error into the confidence interval formula results in  $(2.27, 3.29)$ . We are 95% confident that the average number of days per week that all students from the 2013 YRBSS lifted weights was between 2.27 and 3.29 days.

Because the average of all students from the 2011 YRBSS survey is 3.09, which falls within the range of plausible values from the confidence interval, we cannot say the null hypothesis is implausible. That is, we fail to reject the null hypothesis,  $H_0$ .

#### TIP: Double negatives can sometimes be used in statistics

In many statistical explanations, we use double negatives. For instance, we might say that the null hypothesis is *not implausible* or we *failed to reject* the null hypothesis. Double negatives are used to communicate that while we are not rejecting a position, we are also not saying it is correct.

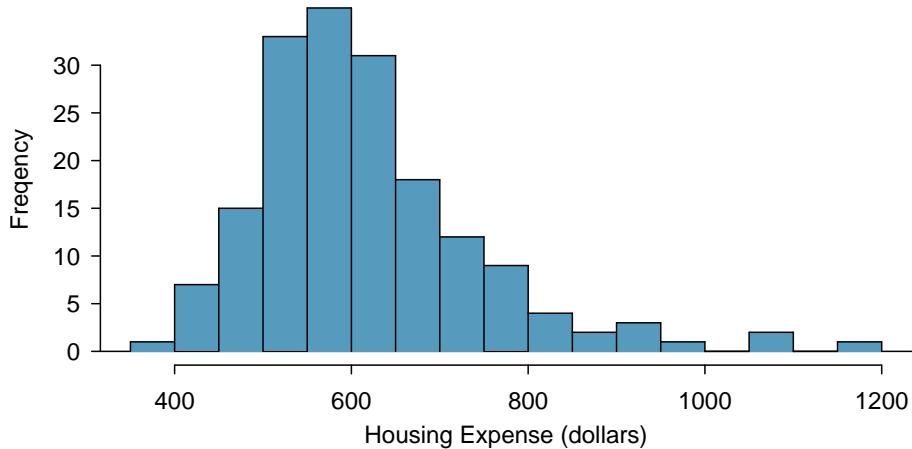


Figure 4.11: Sample distribution of student housing expense. These data are strongly skewed, which we can see by the long right tail with a few notable outliers.

- Ⓐ **Guided Practice 4.20** Colleges frequently provide estimates of student expenses such as housing. A consultant hired by a community college claimed that the average student housing expense was \$650 per month. What are the null and alternative hypotheses to test whether this claim is accurate?<sup>17</sup>
- Ⓑ **Guided Practice 4.21** The community college decides to collect data to evaluate the \$650 per month claim. They take a random sample of 175 students at their school and obtain the data represented in Figure 4.11. Can we apply the normal model to the sample mean?<sup>18</sup>

#### Evaluating the skew condition is challenging

Don't despair if checking the skew condition is difficult or confusing. You aren't alone – nearly all students get frustrated when checking skew. Properly assessing skew takes practice, and you won't be a pro, even at the end of this book.

But this doesn't mean you should give up. Checking skew and the other conditions is extremely important for a responsible data analysis. However, rest assured that evaluating skew isn't something you need to be a master of by the end of the book, though by that time you should be able to properly assess clear cut cases.

<sup>17</sup>  $H_0$ : The average cost is \$650 per month,  $\mu = \$650$ .

$H_A$ : The average cost is different than \$650 per month,  $\mu \neq \$650$ .

<sup>18</sup> Applying the normal model requires that certain conditions are met. Because the data are a simple random sample and the sample (presumably) represents no more than 10% of all students at the college, the observations are independent. The sample size is also sufficiently large ( $n = 175$ ) and the data exhibit strong skew. While the data are strongly skewed, the sample is sufficiently large that this is acceptable, and the normal model may be applied to the sample mean.

- **Example 4.22** The sample mean for student housing is \$616.91 and the sample standard deviation is \$128.65. Construct a 95% confidence interval for the population mean and evaluate the hypotheses of Guided Practice 4.20.

The standard error associated with the mean may be estimated using the sample standard deviation divided by the square root of the sample size. Recall that  $n = 175$  students were sampled.

$$SE = \frac{s}{\sqrt{n}} = \frac{128.65}{\sqrt{175}} = 9.73$$

You showed in Guided Practice 4.21 that the normal model may be applied to the sample mean. This ensures a 95% confidence interval may be accurately constructed:

$$\bar{x} \pm z^* SE \rightarrow 616.91 \pm 1.96 \times 9.73 \rightarrow (597.84, 635.98)$$

Because the null value \$650 is not in the confidence interval, a true mean of \$650 is implausible and we reject the null hypothesis. The data provide statistically significant evidence that the actual average housing expense is less than \$650 per month.

### 4.3.3 Decision errors

Hypothesis tests are not flawless, since we can make a wrong decision in statistical hypothesis tests based on the data. For example, in the court system innocent people are sometimes wrongly convicted and the guilty sometimes walk free. However, the difference is that in statistical hypothesis tests, we have the tools necessary to quantify how often we make such errors.

There are two competing hypotheses: the null and the alternative. In a hypothesis test, we make a statement about which one might be true, but we might choose incorrectly. There are four possible scenarios, which are summarized in Table 4.12.

		Test conclusion	
		do not reject $H_0$	reject $H_0$ in favor of $H_A$
Truth	$H_0$ true	okay	Type 1 Error
	$H_A$ true	Type 2 Error	okay

Table 4.12: Four different scenarios for hypothesis tests.

A **Type 1 Error** is rejecting the null hypothesis when  $H_0$  is actually true. A **Type 2 Error** is failing to reject the null hypothesis when the alternative is actually true.

- **Guided Practice 4.23** In a US court, the defendant is either innocent ( $H_0$ ) or guilty ( $H_A$ ). What does a Type 1 Error represent in this context? What does a Type 2 Error represent? Table 4.12 may be useful.<sup>19</sup>

- **Guided Practice 4.24** How could we reduce the Type 1 Error rate in US courts? What influence would this have on the Type 2 Error rate?<sup>20</sup>

<sup>19</sup>If the court makes a Type 1 Error, this means the defendant is innocent ( $H_0$  true) but wrongly convicted. A Type 2 Error means the court failed to reject  $H_0$  (i.e. failed to convict the person) when she was in fact guilty ( $H_A$  true).

<sup>20</sup>To lower the Type 1 Error rate, we might raise our standard for conviction from “beyond a reasonable doubt” to “beyond a conceivable doubt” so fewer people would be wrongly convicted. However, this would also make it more difficult to convict the people who are actually guilty, so we would make more Type 2 Errors.

- Ⓐ **Guided Practice 4.25** How could we reduce the Type 2 Error rate in US courts? What influence would this have on the Type 1 Error rate?<sup>21</sup>

Exercises 4.23-4.25 provide an important lesson: if we reduce how often we make one type of error, we generally make more of the other type.

Hypothesis testing is built around rejecting or failing to reject the null hypothesis. That is, we do not reject  $H_0$  unless we have strong evidence. But what precisely does *strong evidence* mean? As a general rule of thumb, for those cases where the null hypothesis is actually true, we do not want to incorrectly reject  $H_0$  more than 5% of the time. This corresponds to a **significance level** of 0.05. We often write the significance level using  $\alpha$  (the Greek letter *alpha*):  $\alpha = 0.05$ . We discuss the appropriateness of different significance levels in Section 4.3.6.

If we use a 95% confidence interval to evaluate a hypothesis test where the null hypothesis is true, we will make an error whenever the point estimate is at least 1.96 standard errors away from the population parameter. This happens about 5% of the time (2.5% in each tail). Similarly, using a 99% confidence interval to evaluate a hypothesis is equivalent to a significance level of  $\alpha = 0.01$ .

A confidence interval is, in one sense, simplistic in the world of hypothesis tests. Consider the following two scenarios:

- The null value (the parameter value under the null hypothesis) is in the 95% confidence interval but just barely, so we would not reject  $H_0$ . However, we might like to somehow say, quantitatively, that it was a close decision.
- The null value is very far outside of the interval, so we reject  $H_0$ . However, we want to communicate that, not only did we reject the null hypothesis, but it wasn't even close. Such a case is depicted in Figure 4.13.

In Section 4.3.4, we introduce a tool called the *p-value* that will be helpful in these cases. The p-value method also extends to hypothesis tests where confidence intervals cannot be easily constructed or applied.

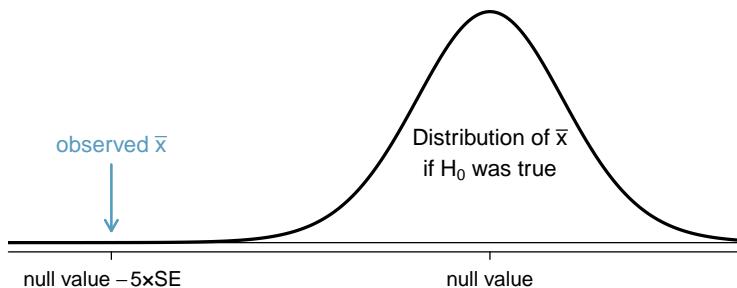


Figure 4.13: It would be helpful to quantify the strength of the evidence against the null hypothesis. In this case, the evidence is extremely strong.

<sup>21</sup>To lower the Type 2 Error rate, we want to convict more guilty people. We could lower the standards for conviction from “beyond a reasonable doubt” to “beyond a little doubt”. Lowering the bar for guilt will also result in more wrongful convictions, raising the Type 1 Error rate.

### 4.3.4 Formal testing using p-values

The p-value is a way of quantifying the strength of the evidence against the null hypothesis and in favor of the alternative. Formally the *p-value* is a conditional probability.

#### p-value

The **p-value** is the probability of observing data at least as favorable to the alternative hypothesis as our current data set, if the null hypothesis is true. We typically use a summary statistic of the data, in this chapter the sample mean, to help compute the p-value and evaluate the hypotheses.

- **Guided Practice 4.26** A poll by the National Sleep Foundation found that college students average about 7 hours of sleep per night. Researchers at a rural school are interested in showing that students at their school sleep longer than seven hours on average, and they would like to demonstrate this using a sample of students. What would be an appropriate skeptical position for this research?<sup>22</sup>

We can set up the null hypothesis for this test as a skeptical perspective: the students at this school average 7 hours of sleep per night. The alternative hypothesis takes a new form reflecting the interests of the research: the students average more than 7 hours of sleep. We can write these hypotheses as

$$H_0: \mu = 7.$$

$$H_A: \mu > 7.$$

Using  $\mu > 7$  as the alternative is an example of a **one-sided** hypothesis test. In this investigation, there is no apparent interest in learning whether the mean is less than 7 hours.<sup>23</sup> Earlier we encountered a **two-sided** hypothesis where we looked for any clear difference, greater than or less than the null value.

Always use a two-sided test unless it was made clear prior to data collection that the test should be one-sided. Switching a two-sided test to a one-sided test after observing the data is dangerous because it can inflate the Type 1 Error rate.

#### TIP: One-sided and two-sided tests

When you are interested in checking for an increase or a decrease, but not both, use a one-sided test. When you are interested in any difference from the null value – an increase or decrease – then the test should be two-sided.

#### TIP: Always write the null hypothesis as an equality

We will find it most useful if we always list the null hypothesis as an equality (e.g.  $\mu = 7$ ) while the alternative always uses an inequality (e.g.  $\mu \neq 7$ ,  $\mu > 7$ , or  $\mu < 7$ ).

<sup>22</sup>A skeptic would have no reason to believe that sleep patterns at this school are different than the sleep patterns at another school.

<sup>23</sup>This is entirely based on the interests of the researchers. Had they been only interested in the opposite case – showing that their students were actually averaging fewer than seven hours of sleep but not interested in showing more than 7 hours – then our setup would have set the alternative as  $\mu < 7$ .

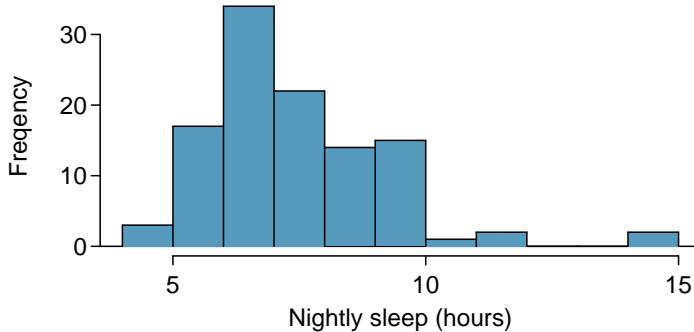


Figure 4.14: Distribution of a night of sleep for 110 college students. These data are strongly skewed.

The researchers at the rural school conducted a simple random sample of  $n = 110$  students on campus. They found that these students averaged 7.42 hours of sleep and the standard deviation of the amount of sleep for the students was 1.75 hours. A histogram of the sample is shown in Figure 4.14.

Before we can use a normal model for the sample mean or compute the standard error of the sample mean, we must verify conditions. (1) Because this is a simple random sample from less than 10% of the student body, the observations are independent. (2) The sample size in the sleep study is sufficiently large since it is greater than 30. (3) The data show strong skew in Figure 4.14 and the presence of a couple of outliers. This skew and the outliers are acceptable for a sample size of  $n = 110$ . With these conditions verified, the normal model can be safely applied to  $\bar{x}$  and we can reasonably calculate the standard error.

- Ⓐ **Guided Practice 4.27** In the sleep study, the sample standard deviation was 1.75 hours and the sample size is 110. Calculate the standard error of  $\bar{x}$ .<sup>24</sup>

The hypothesis test for the sleep study will be evaluated using a significance level of  $\alpha = 0.05$ . We want to consider the data under the scenario that the null hypothesis is true. In this case, the sample mean is from a distribution that is nearly normal and has mean 7 and standard deviation of about  $SE_{\bar{x}} = 0.17$ . Such a distribution is shown in Figure 4.15.

The shaded tail in Figure 4.15 represents the chance of observing such a large mean, conditional on the null hypothesis being true. That is, the shaded tail represents the p-value. We shade all means larger than our sample mean,  $\bar{x} = 7.42$ , because they are more favorable to the alternative hypothesis than the observed mean.

We compute the p-value by finding the tail area of this normal distribution, which we learned to do in Section 3.1. First compute the Z-score of the sample mean,  $\bar{x} = 7.42$ :

$$Z = \frac{\bar{x} - \text{null value}}{SE_{\bar{x}}} = \frac{7.42 - 7}{0.17} = 2.47$$

Using the normal probability table, the lower unshaded area is found to be 0.993. Thus the shaded area is  $1 - 0.993 = 0.007$ . *If the null hypothesis is true, the probability of observing a sample mean at least as large as 7.42 hours for a sample of 110 students is only 0.007.* That is, if the null hypothesis is true, we would not often see such a large mean.

<sup>24</sup>The standard error can be estimated from the sample standard deviation and the sample size:  $SE_{\bar{x}} = \frac{s_{\bar{x}}}{\sqrt{n}} = \frac{1.75}{\sqrt{110}} = 0.17$ .

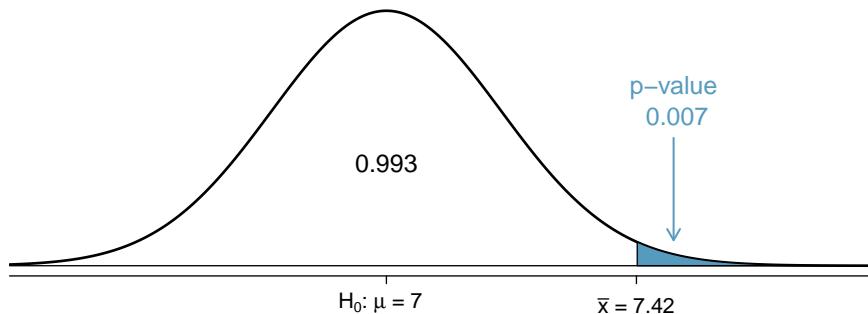


Figure 4.15: If the null hypothesis is true, then the sample mean  $\bar{x}$  came from this nearly normal distribution. The right tail describes the probability of observing such a large sample mean if the null hypothesis is true.

We evaluate the hypotheses by comparing the p-value to the significance level. Because the p-value is less than the significance level ( $p\text{-value} = 0.007 < 0.05 = \alpha$ ), we reject the null hypothesis. What we observed is so unusual with respect to the null hypothesis that it casts serious doubt on  $H_0$  and provides strong evidence favoring  $H_A$ .

#### p-value as a tool in hypothesis testing

The smaller the p-value, the stronger the data favor  $H_A$  over  $H_0$ . A small p-value (usually  $< 0.05$ ) corresponds to sufficient evidence to reject  $H_0$  in favor of  $H_A$ .

#### TIP: It is useful to first draw a picture to find the p-value

It is useful to draw a picture of the distribution of  $\bar{x}$  as though  $H_0$  was true (i.e.  $\mu$  equals the null value), and shade the region (or regions) of sample means that are at least as favorable to the alternative hypothesis. These shaded regions represent the p-value.

The ideas below review the process of evaluating hypothesis tests with p-values:

- The null hypothesis represents a skeptic's position or a position of no difference. We reject this position only if the evidence strongly favors  $H_A$ .
- A small p-value means that if the null hypothesis is true, there is a low probability of seeing a point estimate at least as extreme as the one we saw. We interpret this as strong evidence in favor of the alternative.
- We reject the null hypothesis if the p-value is smaller than the significance level,  $\alpha$ , which is usually 0.05. Otherwise, we fail to reject  $H_0$ .
- We should always state the conclusion of the hypothesis test in plain language so non-statisticians can also understand the results.

The p-value is constructed in such a way that we can directly compare it to the significance level ( $\alpha$ ) to determine whether or not to reject  $H_0$ . This method ensures that the Type 1 Error rate does not exceed the significance level standard.

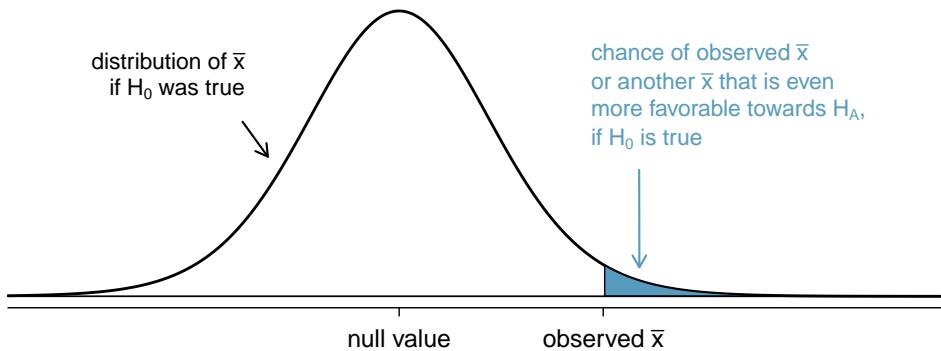


Figure 4.16: To identify the p-value, the distribution of the sample mean is considered as if the null hypothesis was true. Then the p-value is defined and computed as the probability of the observed  $\bar{x}$  or an  $\bar{x}$  even more favorable to  $H_A$  under this distribution.

- Ⓐ **Guided Practice 4.28** If the null hypothesis is true, how often should the p-value be less than 0.05?<sup>25</sup>
- Ⓐ **Guided Practice 4.29** Suppose we had used a significance level of 0.01 in the sleep study. Would the evidence have been strong enough to reject the null hypothesis? (The p-value was 0.007.) What if the significance level was  $\alpha = 0.001$ ?<sup>26</sup>
- Ⓐ **Guided Practice 4.30** Ebay might be interested in showing that buyers on its site tend to pay less than they would for the corresponding new item on Amazon. We'll research this topic for one particular product: a video game called *Mario Kart* for the Nintendo Wii. During early October 2009, Amazon sold this game for \$46.99. Set up an appropriate (one-sided!) hypothesis test to check the claim that Ebay buyers pay less during auctions at this same time.<sup>27</sup>
- Ⓐ **Guided Practice 4.31** During early October 2009, 52 Ebay auctions were recorded for *Mario Kart*.<sup>28</sup> The total prices for the auctions are presented using a histogram in Figure 4.17, and we may like to apply the normal model to the sample mean. Check the three conditions required for applying the normal model: (1) independence, (2) at least 30 observations, and (3) the data are not strongly skewed.<sup>29</sup>

<sup>25</sup>About 5% of the time. If the null hypothesis is true, then the data only has a 5% chance of being in the 5% of data most favorable to  $H_A$ .

<sup>26</sup>We reject the null hypothesis whenever  $p\text{-value} < \alpha$ . Thus, we would still reject the null hypothesis if  $\alpha = 0.01$  but not if the significance level had been  $\alpha = 0.001$ .

<sup>27</sup>The skeptic would say the average is the same on Ebay, and we are interested in showing the average price is lower.

$H_0$ : The average auction price on Ebay is equal to (or more than) the price on Amazon. We write only the equality in the statistical notation:  $\mu_{ebay} = 46.99$ .

$H_A$ : The average price on Ebay is less than the price on Amazon,  $\mu_{ebay} < 46.99$ .

<sup>28</sup>These data were collected by OpenIntro staff.

<sup>29</sup>(1) The independence condition is unclear. *We will make the assumption that the observations are independent, which we should report with any final results.* (2) The sample size is sufficiently large:  $n = 52 \geq 30$ . (3) The data distribution is not strongly skewed; it is approximately symmetric.

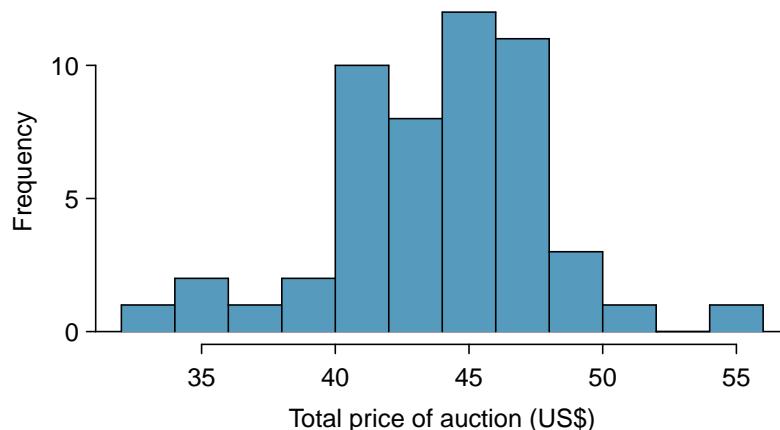


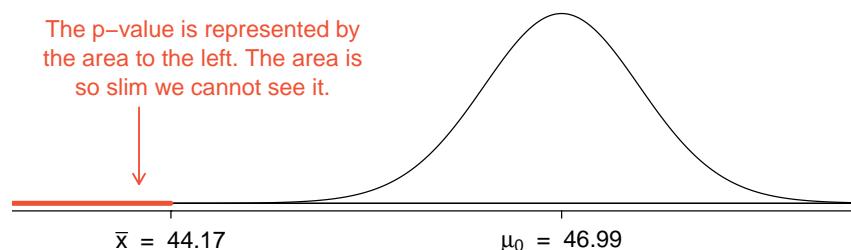
Figure 4.17: A histogram of the total auction prices for 52 Ebay auctions.

- Example 4.32 The average sale price of the 52 Ebay auctions for *Wii Mario Kart* was \$44.17 with a standard deviation of \$4.15. Does this provide sufficient evidence to reject the null hypothesis in Guided Practice 4.30? Use a significance level of  $\alpha = 0.01$ .

---

The hypotheses were set up and the conditions were checked in Exercises 4.30 and 4.31. The next step is to find the standard error of the sample mean and produce a sketch to help find the p-value.

$$SE_{\bar{x}} = s/\sqrt{n} = 4.15/\sqrt{52} = 0.5755$$



Because the alternative hypothesis says we are looking for a smaller mean, we shade the lower tail. We find this shaded area by using the Z-score and normal probability table:  $Z = \frac{44.17 - 46.99}{0.5755} = -4.90$ , which has area less than 0.0002. The area is so small we cannot really see it on the picture. This lower tail area corresponds to the p-value.

Because the p-value is so small – specifically, smaller than  $\alpha = 0.01$  – this provides sufficiently strong evidence to reject the null hypothesis in favor of the alternative. The data provide statistically significant evidence that the average price on Ebay is lower than Amazon's asking price.

### What's so special about 0.05?

It's common to use a threshold of 0.05 to determine whether a result is statistically significant, but why is the most common value 0.05? Maybe the standard significance level should be bigger, or maybe it should be smaller. If you're a little puzzled, that probably means you're reading with a critical eye – good job! We've made a 5-minute task to help clarify *why 0.05*:

[www.openintro.org/why05](http://www.openintro.org/why05)

Sometimes it's also a good idea to deviate from the standard. We'll discuss when to choose a threshold different than 0.05 in Section 4.3.6.

### 4.3.5 Two-sided hypothesis testing with p-values

We now consider how to compute a p-value for a two-sided test. In one-sided tests, we shade the single tail in the direction of the alternative hypothesis. For example, when the alternative had the form  $\mu > 7$ , then the p-value was represented by the upper tail (Figure 4.16). When the alternative was  $\mu < 46.99$ , the p-value was the lower tail (Guided Practice 4.30). In a two-sided test, *we shade two tails* since evidence in either direction is favorable to  $H_A$ .

- Ⓐ **Guided Practice 4.33** Earlier we talked about a research group investigating whether the students at their school slept longer than 7 hours each night. Let's consider a second group of researchers who want to evaluate whether the students at their college differ from the norm of 7 hours. Write the null and alternative hypotheses for this investigation.<sup>30</sup>

- Ⓑ **Example 4.34** The second college randomly samples 122 students and finds a mean of  $\bar{x} = 6.83$  hours and a standard deviation of  $s = 1.8$  hours. Does this provide strong evidence against  $H_0$  in Guided Practice 4.33? Use a significance level of  $\alpha = 0.05$ .

First, we must verify assumptions. (1) A simple random sample of less than 10% of the student body means the observations are independent. (2) The sample size is 122, which is greater than 30. (3) Based on the earlier distribution and what we already know about college student sleep habits, the sample size will be acceptable.

Next we can compute the standard error ( $SE_{\bar{x}} = \frac{s}{\sqrt{n}} = 0.16$ ) of the estimate and create a picture to represent the p-value, shown in Figure 4.18. Both tails are shaded. An estimate of 7.17 or more provides at least as strong of evidence against the null hypothesis and in favor of the alternative as the observed estimate,  $\bar{x} = 6.83$ .

We can calculate the tail areas by first finding the lower tail corresponding to  $\bar{x}$ :

$$Z = \frac{6.83 - 7.00}{0.16} = -1.06 \quad \xrightarrow{\text{table}} \quad \text{left tail} = 0.1446$$

Because the normal model is symmetric, the right tail will have the same area as the left tail. The p-value is found as the sum of the two shaded tails:

$$\text{p-value} = \text{left tail} + \text{right tail} = 2 \times (\text{left tail}) = 0.2892$$

<sup>30</sup>Because the researchers are interested in any difference, they should use a two-sided setup:  $H_0 : \mu = 7$ ,  $H_A : \mu \neq 7$ .

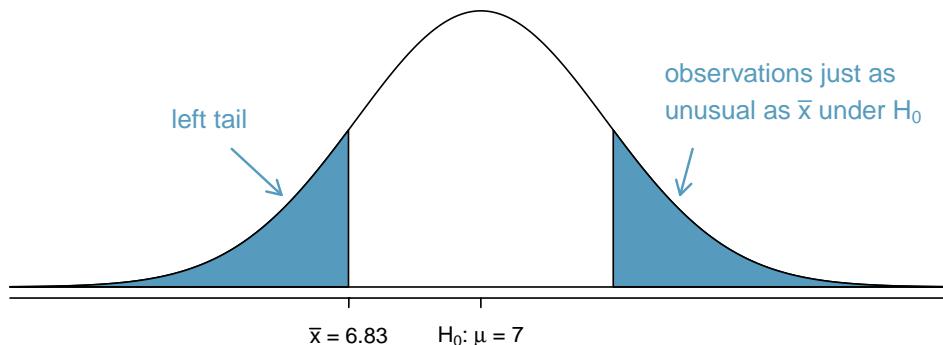


Figure 4.18:  $H_A$  is two-sided, so *both* tails must be counted for the p-value.

This p-value is relatively large (larger than  $\alpha = 0.05$ ), so we should not reject  $H_0$ . That is, if  $H_0$  is true, it would not be very unusual to see a sample mean this far from 7 hours simply due to sampling variation. Thus, we do not have sufficient evidence to conclude that the mean is different than 7 hours.

- **Example 4.35** It is never okay to change two-sided tests to one-sided tests after observing the data. In this example we explore the consequences of ignoring this advice. Using  $\alpha = 0.05$ , we show that freely switching from two-sided tests to one-sided tests will cause us to make twice as many Type 1 Errors as intended.

Suppose the sample mean was larger than the null value,  $\mu_0$  (e.g.  $\mu_0$  would represent 7 if  $H_0: \mu = 7$ ). Then if we can flip to a one-sided test, we would use  $H_A: \mu > \mu_0$ . Now if we obtain any observation with a Z-score greater than 1.65, we would reject  $H_0$ . If the null hypothesis is true, we incorrectly reject the null hypothesis about 5% of the time when the sample mean is above the null value, as shown in Figure 4.19.

Suppose the sample mean was smaller than the null value. Then if we change to a one-sided test, we would use  $H_A: \mu < \mu_0$ . If  $\bar{x}$  had a Z-score smaller than -1.65, we would reject  $H_0$ . If the null hypothesis is true, then we would observe such a case about 5% of the time.

By examining these two scenarios, we can determine that we will make a Type 1 Error  $5\% + 5\% = 10\%$  of the time if we are allowed to swap to the “best” one-sided test for the data. This is twice the error rate we prescribed with our significance level:  $\alpha = 0.05$  (!).

**Caution: One-sided hypotheses are allowed only *before* seeing data**

After observing data, it is tempting to turn a two-sided test into a one-sided test. Avoid this temptation. Hypotheses must be set up *before* observing the data. If they are not, the test should be two-sided.

### 4.3.6 Choosing a significance level

Choosing a significance level for a test is important in many contexts, and the traditional level is 0.05. However, it is often helpful to adjust the significance level based on the

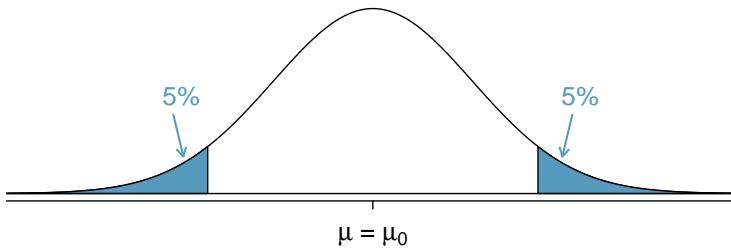


Figure 4.19: The shaded regions represent areas where we would reject  $H_0$  under the bad practices considered in Example 4.35 when  $\alpha = 0.05$ .

application. We may select a level that is smaller or larger than 0.05 depending on the consequences of any conclusions reached from the test.

If making a Type 1 Error is dangerous or especially costly, we should choose a small significance level (e.g. 0.01). Under this scenario we want to be very cautious about rejecting the null hypothesis, so we demand very strong evidence favoring  $H_A$  before we would reject  $H_0$ .

If a Type 2 Error is relatively more dangerous or much more costly than a Type 1 Error, then we should choose a higher significance level (e.g. 0.10). Here we want to be cautious about failing to reject  $H_0$  when the null is actually false.

#### Significance levels should reflect consequences of errors

The significance level selected for a test should reflect the consequences associated with Type 1 and Type 2 Errors.

- **Example 4.36** A car manufacturer is considering a higher quality but more expensive supplier for window parts in its vehicles. They sample a number of parts from their current supplier and also parts from the new supplier. They decide that if the high quality parts will last more than 12% longer, it makes financial sense to switch to this more expensive supplier. Is there good reason to modify the significance level in such a hypothesis test?

The null hypothesis is that the more expensive parts last no more than 12% longer while the alternative is that they do last more than 12% longer. This decision is just one of the many regular factors that have a marginal impact on the car and company. A significance level of 0.05 seems reasonable since neither a Type 1 or Type 2 Error should be dangerous or (relatively) much more expensive.

- **Example 4.37** The same car manufacturer is considering a slightly more expensive supplier for parts related to safety, not windows. If the durability of these safety components is shown to be better than the current supplier, they will switch manufacturers. Is there good reason to modify the significance level in such an evaluation?

The null hypothesis would be that the suppliers' parts are equally reliable. Because safety is involved, the car company should be eager to switch to the slightly more expensive manufacturer (reject  $H_0$ ) even if the evidence of increased safety is only moderately strong. A slightly larger significance level, such as  $\alpha = 0.10$ , might be appropriate.

- Ⓐ **Guided Practice 4.38** A part inside of a machine is very expensive to replace. However, the machine usually functions properly even if this part is broken, so the part is replaced only if we are extremely certain it is broken based on a series of measurements. Identify appropriate hypotheses for this test (in plain language) and suggest an appropriate significance level.<sup>31</sup>

## 4.4 Examining the Central Limit Theorem



The normal model for the sample mean tends to be very good when the sample consists of at least 30 independent observations and the population data are not strongly skewed. The Central Limit Theorem provides the theory that allows us to make this assumption.

### Central Limit Theorem, informal definition

The distribution of  $\bar{x}$  is approximately normal. The approximation can be poor if the sample size is small, but it improves with larger sample sizes.

The Central Limit Theorem states that when the sample size is small, the normal approximation may not be very good. However, as the sample size becomes large, the normal approximation improves. We will investigate three cases to see roughly when the approximation is reasonable.

We consider three data sets: one from a *uniform* distribution, one from an *exponential* distribution, and the other from a *log-normal* distribution. These distributions are shown in the top panels of Figure 4.20. The uniform distribution is symmetric, the exponential distribution may be considered as having moderate skew since its right tail is relatively short (few outliers), and the log-normal distribution is strongly skewed and will tend to produce more apparent outliers.

The left panel in the  $n = 2$  row represents the sampling distribution of  $\bar{x}$  if it is the sample mean of two observations from the uniform distribution shown. The dashed line represents the closest approximation of the normal distribution. Similarly, the center and right panels of the  $n = 2$  row represent the respective distributions of  $\bar{x}$  for data from exponential and log-normal distributions.

- Ⓐ **Guided Practice 4.39** Examine the distributions in each row of Figure 4.20. What do you notice about the normal approximation for each sampling distribution as the sample size becomes larger?<sup>32</sup>

- **Example 4.40** Would the normal approximation be good in all applications where the sample size is at least 30?

Not necessarily. For example, the normal approximation for the log-normal example is questionable for a sample size of 30. Generally, the more skewed a population distribution or the more common the frequency of outliers, the larger the sample required to guarantee the distribution of the sample mean is nearly normal.

<sup>31</sup>Here the null hypothesis is that the part is not broken, and the alternative is that it is broken. If we don't have sufficient evidence to reject  $H_0$ , we would not replace the part. It sounds like failing to fix the part if it is broken ( $H_0$  false,  $H_A$  true) is not very problematic, and replacing the part is expensive. Thus, we should require very strong evidence against  $H_0$  before we replace the part. Choose a small significance level, such as  $\alpha = 0.01$ .

<sup>32</sup>The normal approximation becomes better as larger samples are used.

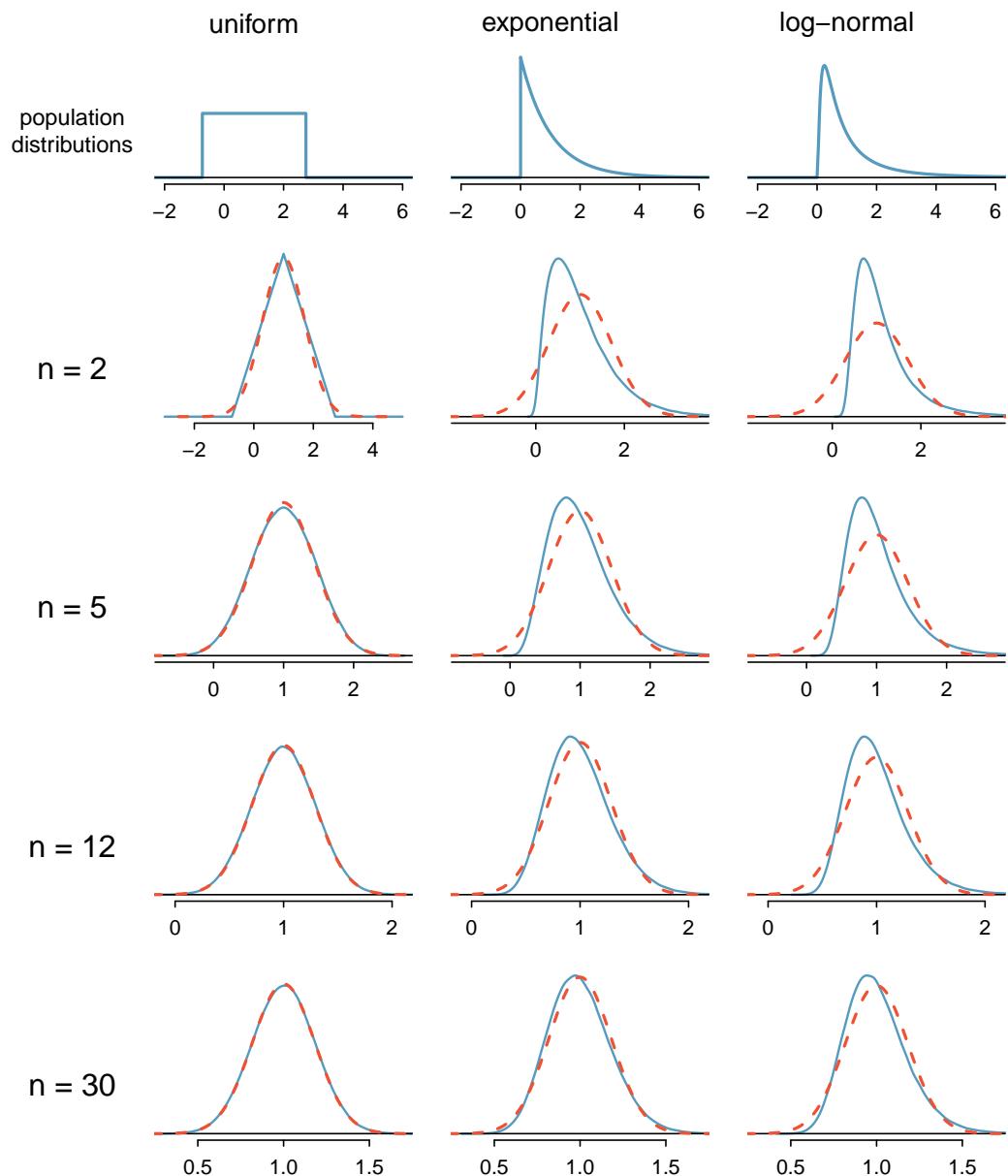


Figure 4.20: Sampling distributions for the mean at different sample sizes and for three different distributions. The dashed red lines show normal distributions.

**TIP: With larger  $n$ , the sampling distribution of  $\bar{x}$  becomes more normal**

As the sample size increases, the normal model for  $\bar{x}$  becomes more reasonable. We can also relax our condition on skew when the sample size is very large.

We discussed in Section 4.1.3 that the sample standard deviation,  $s$ , could be used as a substitute of the population standard deviation,  $\sigma$ , when computing the standard error. This estimate tends to be reasonable when  $n \geq 30$ . We will encounter alternative distributions for smaller sample sizes in Chapters 5 and ??.

- **Example 4.41** Figure 4.21 shows a histogram of 50 observations. These represent winnings and losses from 50 consecutive days of a professional poker player. Can the normal approximation be applied to the sample mean, 90.69?

We should consider each of the required conditions.

- (1) These are referred to as **time series data**, because the data arrived in a particular sequence. If the player wins on one day, it may influence how she plays the next. To make the assumption of independence we should perform careful checks on such data. While the supporting analysis is not shown, no evidence was found to indicate the observations are not independent.
- (2) The sample size is 50, satisfying the sample size condition.
- (3) There are two outliers, one very extreme, which suggests the data are very strongly skewed or very distant outliers may be common for this type of data. Outliers can play an important role and affect the distribution of the sample mean and the estimate of the standard error.

Since we should be skeptical of the independence of observations and the very extreme upper outlier poses a challenge, we should not use the normal model for the sample mean of these 50 observations. If we can obtain a much larger sample, perhaps several hundred observations, then the concerns about skew and outliers would no longer apply.

**Caution: Examine data structure when considering independence**

Some data sets are collected in such a way that they have a natural underlying structure between observations, e.g. when observations occur consecutively. Be especially cautious about independence assumptions regarding such data sets.

**Caution: Watch out for strong skew and outliers**

Strong skew is often identified by the presence of clear outliers. If a data set has prominent outliers, or such observations are somewhat common for the type of data under study, then it is useful to collect a sample with many more than 30 observations if the normal model will be used for  $\bar{x}$ .

You won't be a pro at assessing skew by the end of this book, so just use your best judgement and continue learning. As you develop your statistics skills and encounter tough situations, also consider learning about better ways to analyze skewed data, such as the studentized bootstrap (bootstrap-t), or consult a more experienced statistician.

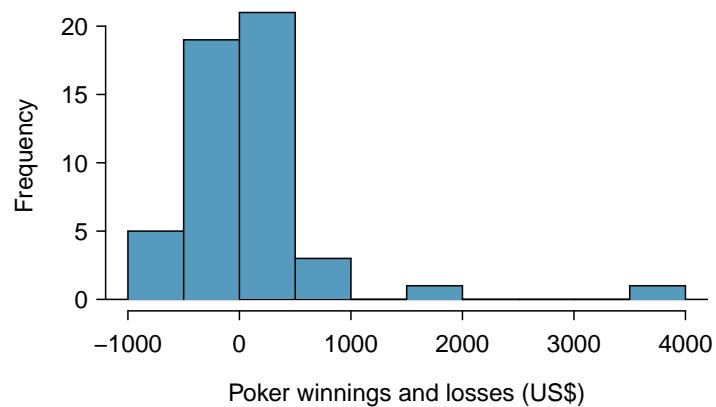


Figure 4.21: Sample distribution of poker winnings. These data include some very clear outliers. These are problematic when considering the normality of the sample mean. For example, outliers are often an indicator of very strong skew.

## 4.5 Exercises

### 4.5.1 Variability in estimates

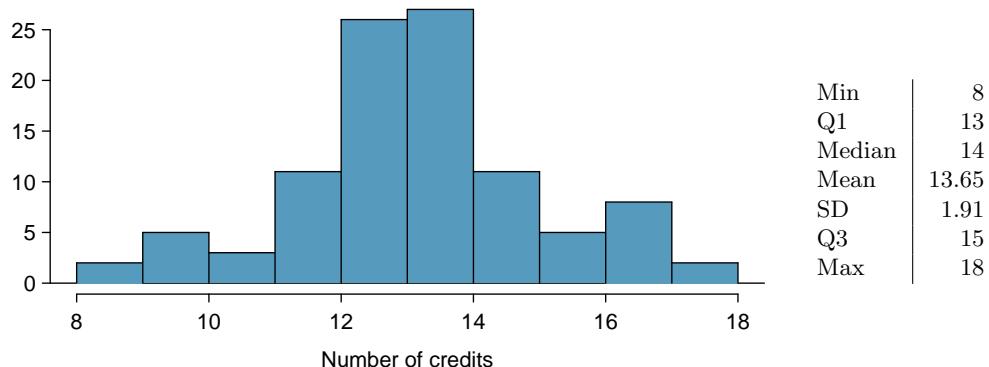
**4.1 Identify the parameter, Part I.** For each of the following situations, state whether the parameter of interest is a mean or a proportion. It may be helpful to examine whether individual responses are numerical or categorical.

- (a) In a survey, one hundred college students are asked how many hours per week they spend on the Internet.
- (b) In a survey, one hundred college students are asked: “What percentage of the time you spend on the Internet is part of your course work?”
- (c) In a survey, one hundred college students are asked whether or not they cited information from Wikipedia in their papers.
- (d) In a survey, one hundred college students are asked what percentage of their total weekly spending is on alcoholic beverages.
- (e) In a sample of one hundred recent college graduates, it is found that 85 percent expect to get a job within one year of their graduation date.

**4.2 Identify the parameter, Part II.** For each of the following situations, state whether the parameter of interest is a mean or a proportion.

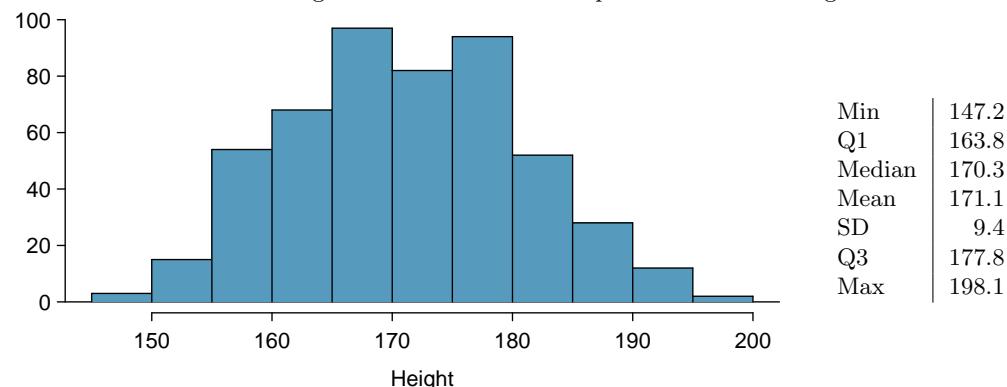
- (a) A poll shows that 64% of Americans personally worry a great deal about federal spending and the budget deficit.
- (b) A survey reports that local TV news has shown a 17% increase in revenue between 2009 and 2011 while newspaper revenues decreased by 6.4% during this time period.
- (c) In a survey, high school and college students are asked whether or not they use geolocation services on their smart phones.
- (d) In a survey, smart phone users are asked whether or not they use a web-based taxi service.
- (e) In a survey, smart phone users are asked how many times they used a web-based taxi service over the last year.

**4.3 College credits.** A college counselor is interested in estimating how many credits a student typically enrolls in each semester. The counselor decides to randomly sample 100 students by using the registrar's database of students. The histogram below shows the distribution of the number of credits taken by these students. Sample statistics for this distribution are also provided.



- What is the point estimate for the average number of credits taken per semester by students at this college? What about the median?
- What is the point estimate for the standard deviation of the number of credits taken per semester by students at this college? What about the IQR?
- Is a load of 16 credits unusually high for this college? What about 18 credits? Explain your reasoning. *Hint:* Observations farther than two standard deviations from the mean are usually considered to be unusual.
- The college counselor takes another random sample of 100 students and this time finds a sample mean of 14.02 units. Should she be surprised that this sample statistic is slightly different than the one from the original sample? Explain your reasoning.
- The sample means given above are point estimates for the mean number of credits taken by all students at that college. What measures do we use to quantify the variability of this estimate (*Hint:* recall that  $SD_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$ )? Compute this quantity using the data from the original sample.

**4.4 Heights of adults.** Researchers studying anthropometry collected body girth measurements and skeletal diameter measurements, as well as age, weight, height and gender, for 507 physically active individuals. The histogram below shows the sample distribution of heights in centimeters.<sup>33</sup>



- What is the point estimate for the average height of active individuals? What about the median? (See the next page for parts (b)-(e).)

<sup>33</sup>G. Heinz et al. "Exploring relationships in body dimensions". In: *Journal of Statistics Education* 11.2 (2003).

- (b) What is the point estimate for the standard deviation of the heights of active individuals? What about the IQR?
- (c) Is a person who is 1m 80cm (180 cm) tall considered unusually tall? And is a person who is 1m 55cm (155cm) considered unusually short? Explain your reasoning.
- (d) The researchers take another random sample of physically active individuals. Would you expect the mean and the standard deviation of this new sample to be the ones given above? Explain your reasoning.
- (e) The sample means obtained are point estimates for the mean height of all active individuals, if the sample of individuals is equivalent to a simple random sample. What measure do we use to quantify the variability of such an estimate (Hint: recall that  $SD_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$ )? Compute this quantity using the data from the original sample under the condition that the data are a simple random sample.

**4.5 Hen eggs.** The distribution of the number of eggs laid by a certain species of hen during their breeding period has a mean of 35 eggs with a standard deviation of 18.2. Suppose a group of researchers randomly samples 45 hens of this species, counts the number of eggs laid during their breeding period, and records the sample mean. They repeat this 1,000 times, and build a distribution of sample means.

- (a) What is this distribution called?
- (b) Would you expect the shape of this distribution to be symmetric, right skewed, or left skewed? Explain your reasoning.
- (c) Calculate the variability of this distribution and state the appropriate term used to refer to this value.
- (d) Suppose the researchers' budget is reduced and they are only able to collect random samples of 10 hens. The sample mean of the number of eggs is recorded, and we repeat this 1,000 times, and build a new distribution of sample means. How will the variability of this new distribution compare to the variability of the original distribution?

**4.6 Art after school.** Elijah and Tyler, two high school juniors, conducted a survey on 15 students at their school, asking the students whether they would like the school to offer an after-school art program, counted the number of "yes" answers, and recorded the sample proportion. 14 out of the 15 students responded "yes". They repeated this 100 times and built a distribution of sample means. (Note that this question requires having reviewed Section ?? on the normal approximation to the binomial distribution.)

- (a) What is this distribution called?
- (b) Would you expect the shape of this distribution to be symmetric, right skewed, or left skewed? Explain your reasoning.
- (c) Calculate the variability of this distribution and state the appropriate term used to refer to this value.
- (d) Suppose that the students were able to recruit a few more friends to help them with sampling, and are now able to collect data from random samples of 25 students. Once again, they record the number of "yes" answers, and record the sample proportion, and repeat this 100 times to build a new distribution of sample proportions. How will the variability of this new distribution compare to the variability of the original distribution?

### 4.5.2 Confidence intervals

**4.7 Chronic illness, Part I.** In 2013, the Pew Research Foundation reported that “45% of U.S. adults report that they live with one or more chronic conditions”.<sup>34</sup> However, this value was based on a sample, so it may not be a perfect estimate for the population parameter of interest on its own. The study reported a standard error of about 1.2%, and a normal model may reasonably be used in this setting. Create a 95% confidence interval for the proportion of U.S. adults who live with one or more chronic conditions. Also interpret the confidence interval in the context of the study.

**4.8 Twitter users and news, Part I.** A poll conducted in 2013 found that 52% of U.S. adult Twitter users get at least some news on Twitter.<sup>35</sup> The standard error for this estimate was 2.4%, and a normal distribution may be used to model the sample proportion. Construct a 99% confidence interval for the fraction of U.S. adult Twitter users who get some news on Twitter, and interpret the confidence interval in context.

**4.9 Chronic illness, Part II.** In 2013, the Pew Research Foundation reported that “45% of U.S. adults report that they live with one or more chronic conditions”, and the standard error for this estimate is 1.2%. Identify each of the following statements as true or false. Provide an explanation to justify each of your answers.

- (a) We can say with certainty that the confidence interval from Exercise 4.7 contains the true percentage of U.S. adults who suffer from a chronic illness.
- (b) If we repeated this study 1,000 times and constructed a 95% confidence interval for each study, then approximately 950 of those confidence intervals would contain the true fraction of U.S. adults who suffer from chronic illnesses.
- (c) The poll provides statistically significant evidence (at the  $\alpha = 0.05$  level) that the percentage of U.S. adults who suffer from chronic illnesses is below 50%.
- (d) Since the standard error is 1.2%, only 1.2% of people in the study communicated uncertainty about their answer.

**4.10 Twitter users and news, Part II.** A poll conducted in 2013 found that 52% of U.S. adult Twitter users get at least some news on Twitter, and the standard error for this estimate was 2.4%. Identify each of the following statements as true or false. Provide an explanation to justify each of your answers.

- (a) The data provide statistically significant evidence that more than half of U.S. adult Twitter users get some news through Twitter. Use a significance level of  $\alpha = 0.01$ .
- (b) Since the standard error is 2.4%, we can conclude that 97.6% of all U.S. adult Twitter users were included in the study.
- (c) If we want to reduce the standard error of the estimate, we should collect less data.
- (d) If we construct a 90% confidence interval for the percentage of U.S. adults Twitter users who get some news through Twitter, this confidence interval will be wider than a corresponding 99% confidence interval.

---

<sup>34</sup>Pew Research Center, Washington, D.C. The Diagnosis Difference, November 26, 2013.

<sup>35</sup>Pew Research Center, Washington, D.C. Twitter News Consumers: Young, Mobile and Educated, November 4, 2013.

**4.11 Relaxing after work.** The 2010 General Social Survey asked the question: “After an average work day, about how many hours do you have to relax or pursue activities that you enjoy?” to a random sample of 1,155 Americans.<sup>36</sup> A 95% confidence interval for the mean number of hours spent relaxing or pursuing activities they enjoy was (1.38, 1.92).

- (a) Interpret this interval in context of the data.
- (b) Suppose another set of researchers reported a confidence interval with a larger margin of error based on the same sample of 1,155 Americans. How does their confidence level compare to the confidence level of the interval stated above?
- (c) Suppose next year a new survey asking the same question is conducted, and this time the sample size is 2,500. Assuming that the population characteristics, with respect to how much time people spend relaxing after work, have not changed much within a year. How will the margin of error of the 95% confidence interval constructed based on data from the new survey compare to the margin of error of the interval stated above?

**4.12 Mental health.** The 2010 General Social Survey asked the question: “For how many days during the past 30 days was your mental health, which includes stress, depression, and problems with emotions, not good?” Based on responses from 1,151 US residents, the survey reported a 95% confidence interval of 3.40 to 4.24 days in 2010.

- (a) Interpret this interval in context of the data.
- (b) What does “95% confident” mean? Explain in the context of the application.
- (c) Suppose the researchers think a 99% confidence level would be more appropriate for this interval. Will this new interval be smaller or larger than the 95% confidence interval?
- (d) If a new survey were to be done with 500 Americans, would the standard error of the estimate be larger, smaller, or about the same. Assume the standard deviation has remained constant since 2010.

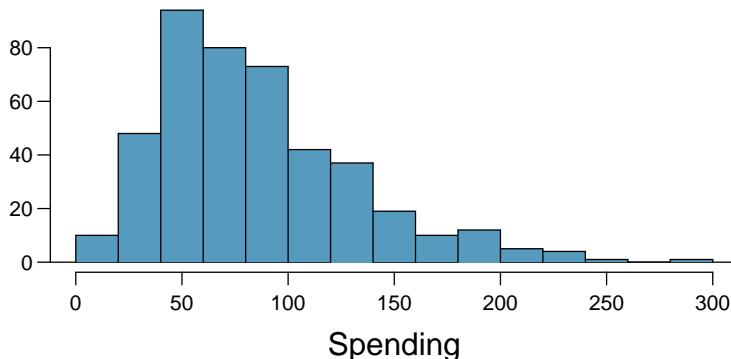
**4.13 Waiting at an ER, Part I.** A hospital administrator hoping to improve wait times decides to estimate the average emergency room waiting time at her hospital. She collects a simple random sample of 64 patients and determines the time (in minutes) between when they checked in to the ER until they were first seen by a doctor. A 95% confidence interval based on this sample is (128 minutes, 147 minutes), which is based on the normal model for the mean. Determine whether the following statements are true or false, and explain your reasoning.

- (a) This confidence interval is not valid since we do not know if the population distribution of the ER wait times is nearly Normal.
- (b) We are 95% confident that the average waiting time of these 64 emergency room patients is between 128 and 147 minutes.
- (c) We are 95% confident that the average waiting time of all patients at this hospital’s emergency room is between 128 and 147 minutes.
- (d) 95% of random samples have a sample mean between 128 and 147 minutes.
- (e) A 99% confidence interval would be narrower than the 95% confidence interval since we need to be more sure of our estimate.
- (f) The margin of error is 9.5 and the sample mean is 137.5.
- (g) In order to decrease the margin of error of a 95% confidence interval to half of what it is now, we would need to double the sample size.

---

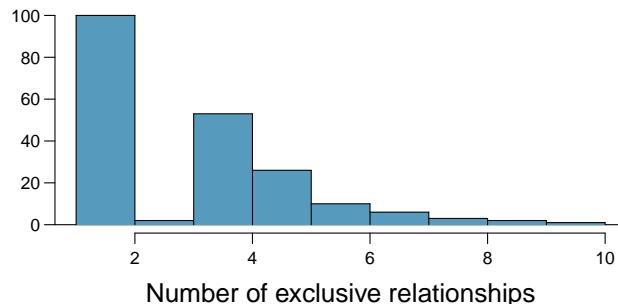
<sup>36</sup>National Opinion Research Center, General Social Survey, 2010.

**4.14 Thanksgiving spending, Part I.** The 2009 holiday retail season, which kicked off on November 27, 2009 (the day after Thanksgiving), had been marked by somewhat lower self-reported consumer spending than was seen during the comparable period in 2008. To get an estimate of consumer spending, 436 randomly sampled American adults were surveyed. Daily consumer spending for the six-day period after Thanksgiving, spanning the Black Friday weekend and Cyber Monday, averaged \$84.71. A 95% confidence interval based on this sample is (\$80.31, \$89.11). Determine whether the following statements are true or false, and explain your reasoning.



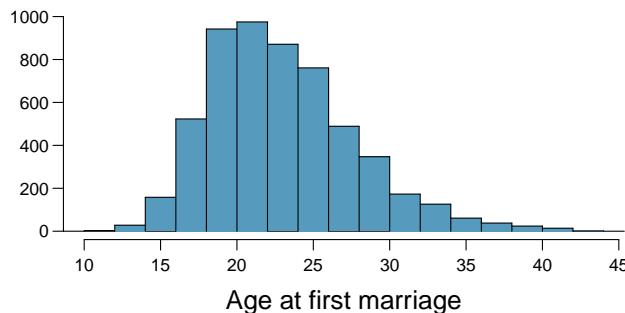
- (a) We are 95% confident that the average spending of these 436 American adults is between \$80.31 and \$89.11.
- (b) This confidence interval is not valid since the distribution of spending in the sample is right skewed.
- (c) 95% of random samples have a sample mean between \$80.31 and \$89.11.
- (d) We are 95% confident that the average spending of all American adults is between \$80.31 and \$89.11.
- (e) A 90% confidence interval would be narrower than the 95% confidence interval since we don't need to be as sure about our estimate.
- (f) In order to decrease the margin of error of a 95% confidence interval to a third of what it is now, we would need to use a sample 3 times larger.
- (g) The margin of error is 4.4.

**4.15 Exclusive relationships.** A survey conducted on a reasonably random sample of 203 undergraduates asked, among many other questions, about the number of exclusive relationships these students have been in. The histogram below shows the distribution of the data from this sample. The sample average is 3.2 with a standard deviation of 1.97.



Estimate the average number of exclusive relationships Duke students have been in using a 90% confidence interval and interpret this interval in context. Check any conditions required for inference, and note any assumptions you must make as you proceed with your calculations and conclusions.

**4.16 Age at first marriage, Part I.** The National Survey of Family Growth conducted by the Centers for Disease Control gathers information on family life, marriage and divorce, pregnancy, infertility, use of contraception, and men's and women's health. One of the variables collected on this survey is the age at first marriage. The histogram below shows the distribution of ages at first marriage of 5,534 randomly sampled women between 2006 and 2010. The average age at first marriage among these women is 23.44 with a standard deviation of 4.72.<sup>37</sup>



Estimate the average age at first marriage of women using a 95% confidence interval, and interpret this interval in context. Discuss any relevant assumptions.

### 4.5.3 Hypothesis testing

**4.17 Identify hypotheses, Part I.** Write the null and alternative hypotheses in words and then symbols for each of the following situations.

- (a) New York is known as “the city that never sleeps”. A random sample of 25 New Yorkers were asked how much sleep they get per night. Do these data provide convincing evidence that New Yorkers on average sleep less than 8 hours a night?
- (b) Employers at a firm are worried about the effect of March Madness, a basketball championship held each spring in the US, on employee productivity. They estimate that on a regular business day employees spend on average 15 minutes of company time checking personal email, making personal phone calls, etc. They also collect data on how much company time employees spend on such non- business activities during March Madness. They want to determine if these data provide convincing evidence that employee productivity decreases during March Madness.

**4.18 Identify hypotheses, Part II.** Write the null and alternative hypotheses in words and using symbols for each of the following situations.

- (a) Since 2008, chain restaurants in California have been required to display calorie counts of each menu item. Prior to menus displaying calorie counts, the average calorie intake of diners at a restaurant was 1100 calories. After calorie counts started to be displayed on menus, a nutritionist collected data on the number of calories consumed at this restaurant from a random sample of diners. Do these data provide convincing evidence of a difference in the average calorie intake of diners at this restaurant?
- (b) Based on the performance of those who took the GRE exam between July 1, 2004 and June 30, 2007, the average Verbal Reasoning score was calculated to be 462. In 2011 the average verbal score was slightly higher. Do these data provide convincing evidence that the average GRE Verbal Reasoning score has changed since 2004?

---

<sup>37</sup>Centers for Disease Control and Prevention, National Survey of Family Growth, 2010.

**4.19 Online communication.** A study suggests that the average college student spends 10 hours per week communicating with others online. You believe that this is an underestimate and decide to collect your own sample for a hypothesis test. You randomly sample 60 students from your dorm and find that on average they spent 13.5 hours a week communicating with others online. A friend of yours, who offers to help you with the hypothesis test, comes up with the following set of hypotheses. Indicate any errors you see.

$$H_0 : \bar{x} < 10 \text{ hours}$$

$$H_A : \bar{x} > 13.5 \text{ hours}$$

**4.20 Age at first marriage, Part II.** Exercise 4.16 presents the results of a 2006 - 2010 survey showing that the average age of women at first marriage is 23.44. Suppose a social scientist believes that this value has increased in 2012, but she would also be interested if she found a decrease. Below is how she set up her hypotheses. Indicate any errors you see.

$$H_0 : \bar{x} = 23.44 \text{ years old}$$

$$H_A : \bar{x} > 23.44 \text{ years old}$$

**4.21 Waiting at an ER, Part II.** Exercise 4.13 provides a 95% confidence interval for the mean waiting time at an emergency room (ER) of (128 minutes, 147 minutes). Answer the following questions based on this interval.

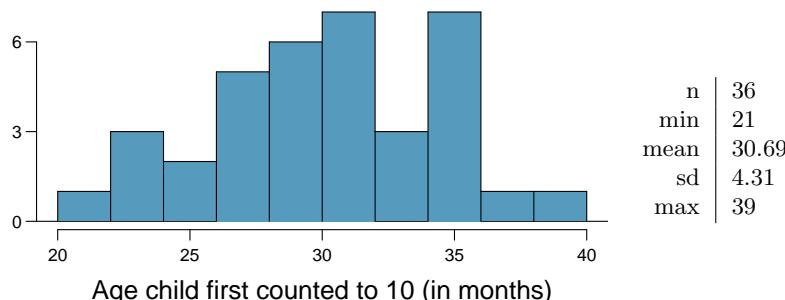
- (a) A local newspaper claims that the average waiting time at this ER exceeds 3 hours. Is this claim supported by the confidence interval? Explain your reasoning.
- (b) The Dean of Medicine at this hospital claims the average wait time is 2.2 hours. Is this claim supported by the confidence interval? Explain your reasoning.
- (c) Without actually calculating the interval, determine if the claim of the Dean from part (b) would be supported based on a 99% confidence interval?

**4.22 Thanksgiving spending, Part II.** Exercise 4.14 provides a 95% confidence interval for the average spending by American adults during the six-day period after Thanksgiving 2009: (\$80.31, \$89.11).

- (a) A local news anchor claims that the average spending during this period in 2009 was \$100. What do you think of her claim?
- (b) Would the news anchor's claim be considered reasonable based on a 90% confidence interval? Why or why not? (Do not actually calculate the interval.)

**4.23 Nutrition labels.** The nutrition label on a bag of potato chips says that a one ounce (28 gram) serving of potato chips has 130 calories and contains ten grams of fat, with three grams of saturated fat. A random sample of 35 bags yielded a sample mean of 134 calories with a standard deviation of 17 calories. Is there evidence that the nutrition label does not provide an accurate measure of calories in the bags of potato chips? We have verified the independence, sample size, and skew conditions are satisfied.

**4.24 Gifted children, Part I.** Researchers investigating characteristics of gifted children collected data from schools in a large city on a random sample of thirty-six children who were identified as gifted children soon after they reached the age of four. The following histogram shows the distribution of the ages (in months) at which these children first counted to 10 successfully. Also provided are some sample statistics.<sup>38</sup>



- Are conditions for inference satisfied?
- Suppose you read online that children first count to 10 successfully when they are 32 months old, on average. Perform a hypothesis test to evaluate if these data provide convincing evidence that the average age at which gifted children first count to 10 successfully is less than the general average of 32 months. Use a significance level of 0.10.
- Interpret the p-value in context of the hypothesis test and the data.
- Calculate a 90% confidence interval for the average age at which gifted children first count to 10 successfully.
- Do your results from the hypothesis test and the confidence interval agree? Explain.

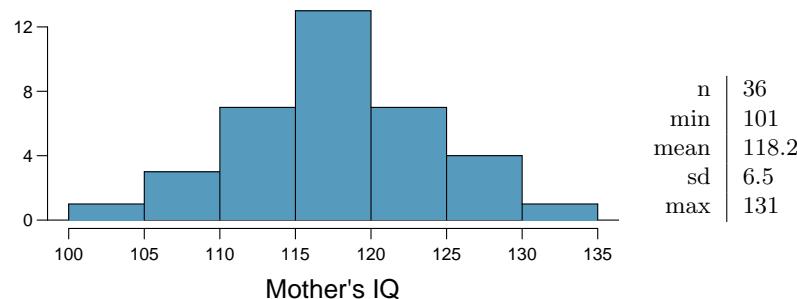
**4.25 Waiting at an ER, Part III.** The hospital administrator mentioned in Exercise 4.13 randomly selected 64 patients and measured the time (in minutes) between when they checked in to the ER and the time they were first seen by a doctor. The average time is 137.5 minutes and the standard deviation is 39 minutes. She is getting grief from her supervisor on the basis that the wait times in the ER has increased greatly from last year's average of 127 minutes. However, she claims that the increase is probably just due to chance.

- Are conditions for inference met? Note any assumptions you must make to proceed.
- Using a significance level of  $\alpha = 0.05$ , is the change in wait times statistically significant? Use a two-sided test since it seems the supervisor had to inspect the data before she suggested an increase occurred.
- Would the conclusion of the hypothesis test change if the significance level was changed to  $\alpha = 0.01$ ?

---

<sup>38</sup>F.A. Graybill and H.K. Iyer. *Regression Analysis: Concepts and Applications*. Duxbury Press, 1994, pp. 511–516.

**4.26 Gifted children, Part II.** Exercise 4.24 describes a study on gifted children. In this study, along with variables on the children, the researchers also collected data on the mother's and father's IQ of the 36 randomly sampled gifted children. The histogram below shows the distribution of mother's IQ. Also provided are some sample statistics.



- (a) Perform a hypothesis test to evaluate if these data provide convincing evidence that the average IQ of mothers of gifted children is different than the average IQ for the population at large, which is 100. Use a significance level of 0.10.
- (b) Calculate a 90% confidence interval for the average IQ of mothers of gifted children.
- (c) Do your results from the hypothesis test and the confidence interval agree? Explain.

**4.27 Working backwards, one-sided.** You are given the following hypotheses:

$$H_0 : \mu = 30$$

$$H_A : \mu > 30$$

We know that the sample standard deviation is 10 and the sample size is 70. For what sample mean would the p-value be equal to 0.05? Assume that all conditions necessary for inference are satisfied.

**4.28 Working backwards, two-sided.** You are given the following hypotheses:

$$H_0 : \mu = 30$$

$$H_A : \mu \neq 30$$

We know that the sample standard deviation is 10 and the sample size is 70. For what sample mean would the p-value be equal to 0.05? Assume that all conditions necessary for inference are satisfied.

**4.29 Testing for Fibromyalgia.** A patient named Diana was diagnosed with Fibromyalgia, a long-term syndrome of body pain, and was prescribed anti-depressants. Being the skeptic that she is, Diana didn't initially believe that anti-depressants would help her symptoms. However after a couple months of being on the medication she decides that the anti-depressants are working, because she feels like her symptoms are in fact getting better.

- (a) Write the hypotheses in words for Diana's skeptical position when she started taking the anti-depressants.
- (b) What is a Type 1 Error in this context?
- (c) What is a Type 2 Error in this context?

**4.30 Testing for food safety.** A food safety inspector is called upon to investigate a restaurant with a few customer reports of poor sanitation practices. The food safety inspector uses a hypothesis testing framework to evaluate whether regulations are not being met. If he decides the restaurant is in gross violation, its license to serve food will be revoked.

- (a) Write the hypotheses in words.
- (b) What is a Type 1 Error in this context?
- (c) What is a Type 2 Error in this context?
- (d) Which error is more problematic for the restaurant owner? Why?
- (e) Which error is more problematic for the diners? Why?
- (f) As a diner, would you prefer that the food safety inspector requires strong evidence or very strong evidence of health concerns before revoking a restaurant's license? Explain your reasoning.

**4.31 Which is higher?** In each part below, there is a value of interest and two scenarios (I and II). For each part, report if the value of interest is larger under scenario I, scenario II, or whether the value is equal under the scenarios.

- (a) The standard error of  $\bar{x}$  when  $s = 120$  and (I)  $n = 25$  or (II)  $n = 125$ .
- (b) The margin of error of a confidence interval when the confidence level is (I) 90% or (II) 80%.
- (c) The p-value for a Z-statistic of 2.5 when (I)  $n = 500$  or (II)  $n = 1000$ .
- (d) The probability of making a Type 2 Error when the alternative hypothesis is true and the significance level is (I) 0.05 or (II) 0.10.

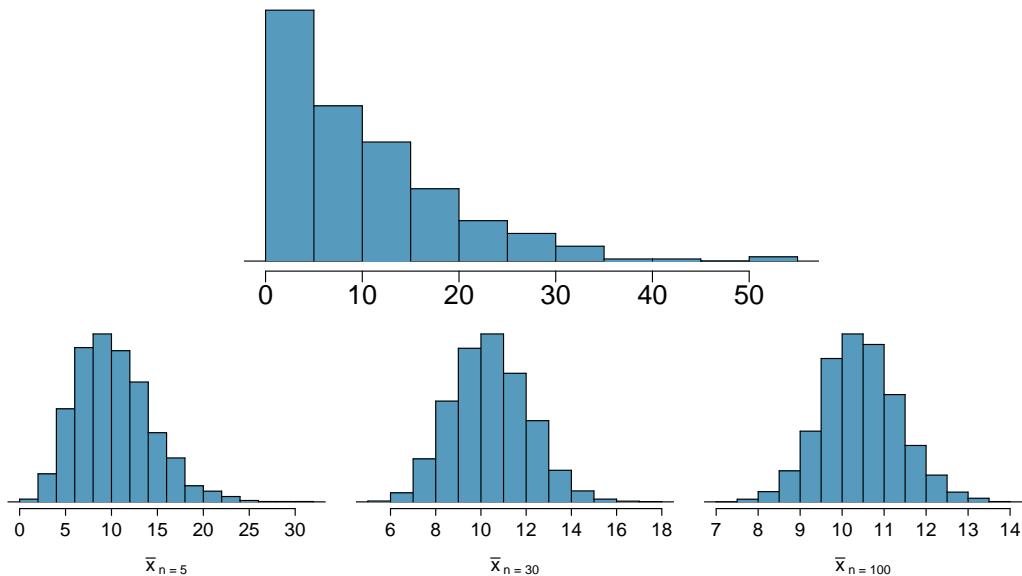
**4.32 True or false.** Determine if the following statements are true or false, and explain your reasoning. If false, state how it could be corrected.

- (a) If a given value (for example, the null hypothesized value of a parameter) is within a 95% confidence interval, it will also be within a 99% confidence interval.
- (b) Decreasing the significance level ( $\alpha$ ) will increase the probability of making a Type 1 Error.
- (c) Suppose the null hypothesis is  $\mu = 5$  and we fail to reject  $H_0$ . Under this scenario, the true population mean is 5.
- (d) If the alternative hypothesis is true, then the probability of making a Type 2 Error and the power of a test add up to 1.
- (e) With large sample sizes, even small differences between the null value and the true value of the parameter, a difference often called the effect size , will be identified as statistically significant.

#### 4.5.4 Examining the Central Limit Theorem

**4.33 Ages of pennies.** The histogram below shows the distribution of ages of pennies at a bank.

- Describe the distribution.
- Sampling distributions for means from simple random samples of 5, 30, and 100 pennies is shown in the histograms below. Describe the shapes of these distributions and comment on whether they look like what you would expect to see based on the Central Limit Theorem.
- The mean age of the pennies is 10.44 years, with a standard deviation of 9.2 years. Using the Central Limit Theorem, calculate the means and standard deviations of the distribution of means from random samples of size 5, 30, and 100. Comment on whether the sampling distributions shown in part (b) agree with the values you compute.



**4.34 CLT.** Define the term “sampling distribution” of the mean, and describe how the shape, center, and spread of the sampling distribution of the mean change as sample size increases.



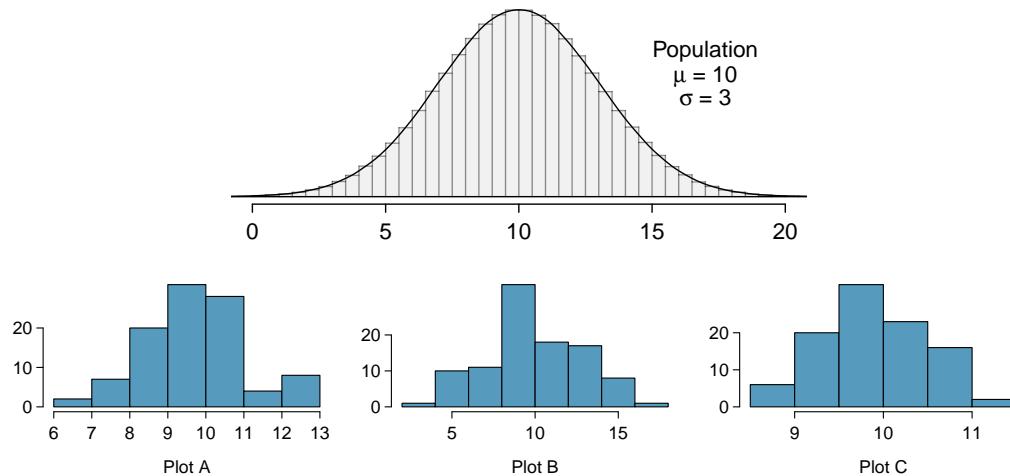
**4.35 Housing prices.** A housing survey was conducted to determine the price of a typical home in Topanga, CA. The mean price of a house was roughly \$1.3 million with a standard deviation of \$300,000. There were no houses listed below \$600,000 but a few houses above \$3 million.

- Is the distribution of housing prices in Topanga symmetric, right skewed, or left skewed? Hint: Sketch the distribution.
- Would you expect most houses in Topanga to cost more or less than \$1.3 million?
- Can we estimate the probability that a randomly chosen house in Topanga costs more than \$1.4 million using the normal distribution?
- What is the probability that the mean of 60 randomly chosen houses in Topanga is more than \$1.4 million?
- How would doubling the sample size affect the standard deviation of the mean?

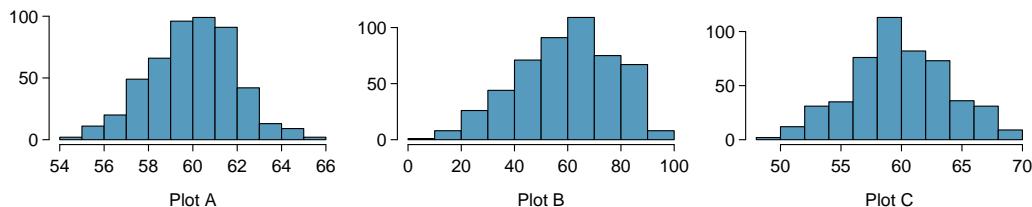
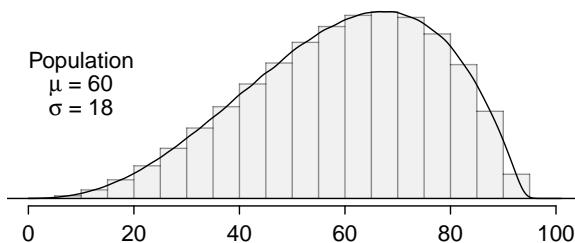
**4.36 Stats final scores.** Each year about 1500 students take the introductory statistics course at a large university. This year scores on the final exam are distributed with a median of 74 points, a mean of 70 points, and a standard deviation of 10 points. There are no students who scored above 100 (the maximum score attainable on the final) but a few students scored below 20 points.

- Is the distribution of scores on this final exam symmetric, right skewed, or left skewed?
- Would you expect most students to have scored above or below 70 points?
- Can we calculate the probability that a randomly chosen student scored above 75 using the normal distribution?
- What is the probability that the average score for a random sample of 40 students is above 75?
- How would cutting the sample size in half affect the standard deviation of the mean?

**4.37 Identify distributions, Part I.** Four plots are presented below. The plot at the top is a distribution for a population. The mean is 10 and the standard deviation is 3. Also shown below is a distribution of (1) a single random sample of 100 values from this population, (2) a distribution of 100 sample means from random samples with size 5, and (3) a distribution of 100 sample means from random samples with size 25. Determine which plot (A, B, or C) is which and explain your reasoning.



**4.38 Identify distributions, Part II.** Four plots are presented below. The plot at the top is a distribution for a population. The mean is 60 and the standard deviation is 18. Also shown below is a distribution of (1) a single random sample of 500 values from this population, (2) a distribution of 500 sample means from random samples of each size 18, and (3) a distribution of 500 sample means from random samples of each size 81. Determine which plot (A, B, or C) is which and explain your reasoning.



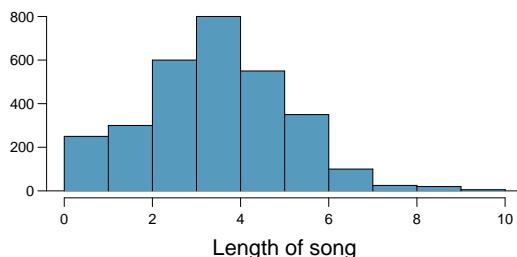
**4.39 Weights of pennies.** The distribution of weights of United States pennies is approximately normal with a mean of 2.5 grams and a standard deviation of 0.03 grams.

- (a) What is the probability that a randomly chosen penny weighs less than 2.4 grams?
- (b) Describe the sampling distribution of the mean weight of 10 randomly chosen pennies.
- (c) What is the probability that the mean weight of 10 pennies is less than 2.4 grams?
- (d) Sketch the two distributions (population and sampling) on the same scale.
- (e) Could you estimate the probabilities from (a) and (c) if the weights of pennies had a skewed distribution?

**4.40 CFLBs.** A manufacturer of compact fluorescent light bulbs advertises that the distribution of the lifespans of these light bulbs is nearly normal with a mean of 9,000 hours and a standard deviation of 1,000 hours.

- (a) What is the probability that a randomly chosen light bulb lasts more than 10,500 hours?
- (b) Describe the distribution of the mean lifespan of 15 light bulbs.
- (c) What is the probability that the mean lifespan of 15 randomly chosen light bulbs is more than 10,500 hours?
- (d) Sketch the two distributions (population and sampling) on the same scale.
- (e) Could you estimate the probabilities from parts (a) and (c) if the lifespans of light bulbs had a skewed distribution?

**4.41 Songs on an iPod.** Suppose an iPod has 3,000 songs. The histogram below shows the distribution of the lengths of these songs. We also know that, for this iPod, the mean length is 3.45 minutes and the standard deviation is 1.63 minutes.



- Calculate the probability that a randomly selected song lasts more than 5 minutes.
- You are about to go for an hour run and you make a random playlist of 15 songs. What is the probability that your playlist lasts for the entire duration of your run? *Hint:* If you want the playlist to last 60 minutes, what should be the minimum average length of a song?
- You are about to take a trip to visit your parents and the drive is 6 hours. You make a random playlist of 100 songs. What is the probability that your playlist lasts the entire drive?

**4.42 Spray paint.** Suppose the area that can be painted using a single can of spray paint is slightly variable and follows a nearly normal distribution with a mean of 25 square feet and a standard deviation of 3 square feet.

- What is the probability that the area covered by a can of spray paint is more than 27 square feet?
- Suppose you want to spray paint an area of 540 square feet using 20 cans of spray paint. On average, how many square feet must each can be able to cover to spray paint all 540 square feet?
- What is the probability that you can cover a 540 square feet area using 20 cans of spray paint?
- If the area covered by a can of spray paint had a slightly skewed distribution, could you still calculate the probabilities in parts (a) and (c) using the normal distribution?

# Chapter 5

## Inference for numerical data

Chapter 4 introduced a framework for statistical inference based on confidence intervals and hypotheses. In this chapter, we encounter several new point estimates and scenarios. In each case, the inference ideas remain the same:

1. Determine which point estimate or test statistic is useful.
2. Identify an appropriate distribution for the point estimate or test statistic.
3. Apply the ideas from Chapter 4 using the distribution from step 2.

### 5.1 One-sample means with the $t$ -distribution



We required a large sample in Chapter 4 for two reasons:

1. The sampling distribution of  $\bar{x}$  tends to be more normal when the sample is large.
2. The calculated standard error is typically very accurate when using a large sample.

So what should we do when the sample size is small? As we'll discuss in Section 5.1.1, if the population data are nearly normal, then  $\bar{x}$  will also follow a normal distribution, which addresses the first problem. The accuracy of the standard error is trickier, and for this challenge we'll introduce a new distribution called the  $t$ -distribution.

While we emphasize the use of the  $t$ -distribution for small samples, this distribution is also generally used for large samples, where it produces similar results to those from the normal distribution.

#### 5.1.1 The normality condition

A special case of the Central Limit Theorem ensures the distribution of sample means will be nearly normal, regardless of sample size, when the data come from a nearly normal distribution.

##### Central Limit Theorem for normal data

The sampling distribution of the mean is nearly normal when the sample observations are independent and come from a nearly normal distribution. This is true for any sample size.

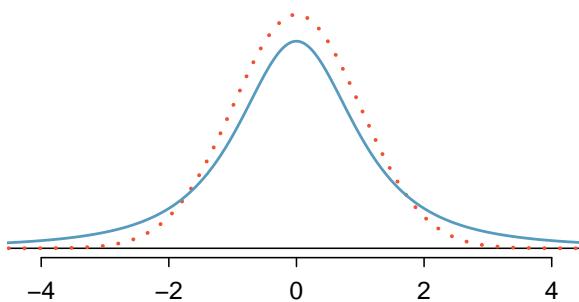


Figure 5.1: Comparison of a  $t$ -distribution (solid line) and a normal distribution (dotted line).

While this seems like a very helpful special case, there is one small problem. It is inherently difficult to verify normality in small data sets.

#### **Caution: Checking the normality condition**

We should exercise caution when verifying the normality condition for small samples. It is important to not only examine the data but also think about where the data come from. For example, ask: would I expect this distribution to be symmetric, and am I confident that outliers are rare?

You may relax the normality condition as the sample size goes up. If the sample size is 10 or more, slight skew is not problematic. Once the sample size hits about 30, then moderate skew is reasonable. Data with strong skew or outliers require a more cautious analysis.

### 5.1.2 Introducing the $t$ -distribution

In the cases where we will use a small sample to calculate the standard error, it will be useful to rely on a new distribution for inference calculations: the  $t$ -distribution. A  $t$ -distribution, shown as a solid line in Figure 5.1, has a bell shape. However, its tails are thicker than the normal model's. This means observations are more likely to fall beyond two standard deviations from the mean than under the normal distribution.<sup>1</sup> While our estimate of the standard error will be a little less accurate when we are analyzing a small data set, these extra thick tails of the  $t$ -distribution are exactly the correction we need to resolve the problem of a poorly estimated standard error.

The  $t$ -distribution, always centered at zero, has a single parameter: degrees of freedom. The **degrees of freedom (df)** describe the precise form of the bell-shaped  $t$ -distribution. Several  $t$ -distributions are shown in Figure 5.2. When there are more degrees of freedom, the  $t$ -distribution looks very much like the standard normal distribution.

---

<sup>1</sup>The standard deviation of the  $t$ -distribution is actually a little more than 1. However, it is useful to always think of the  $t$ -distribution as having a standard deviation of 1 in all of our applications.

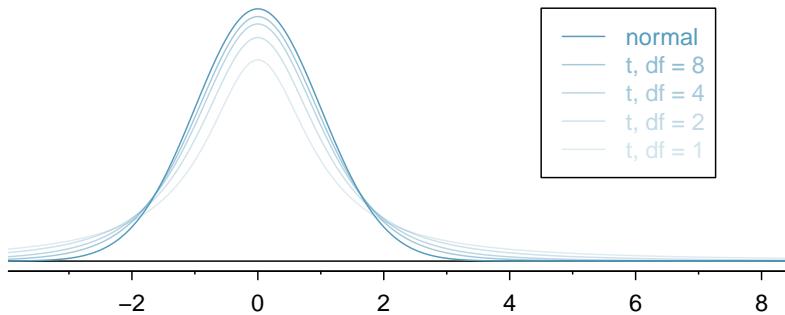


Figure 5.2: The larger the degrees of freedom, the more closely the  $t$ -distribution resembles the standard normal model.

### Degrees of freedom (df)

The degrees of freedom describe the shape of the  $t$ -distribution. The larger the degrees of freedom, the more closely the distribution approximates the normal model.

When the degrees of freedom is about 30 or more, the  $t$ -distribution is nearly indistinguishable from the normal distribution. In Section 5.1.3, we relate degrees of freedom to sample size.

It's very useful to become familiar with the  $t$ -distribution, because it allows us greater flexibility than the normal distribution when analyzing numerical data. We use a ***t*-table**, partially shown in Table 5.3, in place of the normal probability table. A larger  $t$ -table is in Appendix B.2 on page 255. In practice, it's more common to use statistical software instead of a table, and you can see some of these options at

[www.openintro.org/stat/prob-tables](http://www.openintro.org/stat/prob-tables)

Each row in the  $t$ -table represents a  $t$ -distribution with different degrees of freedom. The columns correspond to tail probabilities. For instance, if we know we are working with the  $t$ -distribution with  $df = 18$ , we can examine row 18, which is highlighted in Table 5.3. If we want the value in this row that identifies the cutoff for an upper tail of 10%, we can look in the column where *one tail* is 0.100. This cutoff is 1.33. If we had wanted the cutoff for the lower 10%, we would use -1.33. Just like the normal distribution, all  $t$ -distributions are symmetric.

- **Example 5.1** What proportion of the  $t$ -distribution with 18 degrees of freedom falls below -2.10?

Just like a normal probability problem, we first draw the picture in Figure 5.4 and shade the area below -2.10. To find this area, we identify the appropriate row:  $df = 18$ . Then we identify the column containing the absolute value of -2.10; it is the third column. Because we are looking for just one tail, we examine the top line of the table, which shows that a one tail area for a value in the third row corresponds to 0.025. About 2.5% of the distribution falls below -2.10. In the next example we encounter a case where the exact  $t$  value is not listed in the table.

one tail	0.100	0.050	0.025	0.010	0.005
two tails	0.200	0.100	0.050	0.020	0.010
<i>df</i>					
1	3.08	6.31	12.71	31.82	63.66
2	1.89	2.92	4.30	6.96	9.92
3	1.64	2.35	3.18	4.54	5.84
:	:	:	:	:	:
17	1.33	1.74	2.11	2.57	2.90
<b>18</b>	<b>1.33</b>	<b>1.73</b>	<b>2.10</b>	<b>2.55</b>	<b>2.88</b>
19	1.33	1.73	2.09	2.54	2.86
20	1.33	1.72	2.09	2.53	2.85
:	:	:	:	:	:
400	1.28	1.65	1.97	2.34	2.59
500	1.28	1.65	1.96	2.33	2.59
$\infty$	1.28	1.64	1.96	2.33	2.58

Table 5.3: An abbreviated look at the  $t$ -table. Each row represents a different  $t$ -distribution. The columns describe the cutoffs for specific tail areas. The row with  $df = 18$  has been highlighted.

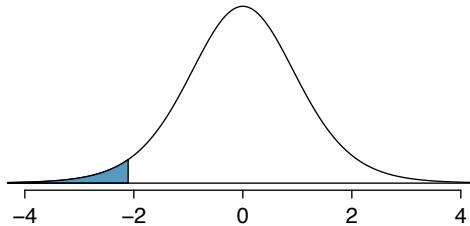


Figure 5.4: The  $t$ -distribution with 18 degrees of freedom. The area below -2.10 has been shaded.

- Example 5.2 A  $t$ -distribution with 20 degrees of freedom is shown in the left panel of Figure 5.5. Estimate the proportion of the distribution falling above 1.65.

We identify the row in the  $t$ -table using the degrees of freedom:  $df = 20$ . Then we look for 1.65; it is not listed. It falls between the first and second columns. Since these values bound 1.65, their tail areas will bound the tail area corresponding to 1.65. We identify the one tail area of the first and second columns, 0.050 and 0.10, and we conclude that between 5% and 10% of the distribution is more than 1.65 standard deviations above the mean. If we like, we can identify the precise area using statistical software: 0.0573.

- Example 5.3 A  $t$ -distribution with 2 degrees of freedom is shown in the right panel of Figure 5.5. Estimate the proportion of the distribution falling more than 3 units from the mean (above or below).

As before, first identify the appropriate row:  $df = 2$ . Next, find the columns that capture 3; because  $2.92 < 3 < 4.30$ , we use the second and third columns. Finally, we find bounds for the tail areas by looking at the two tail values: 0.05 and 0.10. We use the two tail values because we are looking for two (symmetric) tails.

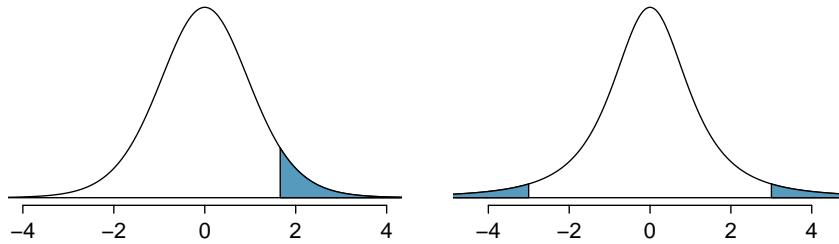


Figure 5.5: Left: The  $t$ -distribution with 20 degrees of freedom, with the area above 1.65 shaded. Right: The  $t$ -distribution with 2 degrees of freedom, with the area further than 3 units from 0 shaded.

Ⓐ **Guided Practice 5.4** What proportion of the  $t$ -distribution with 19 degrees of freedom falls above -1.79 units?<sup>2</sup>

### 5.1.3 Conditions for using the $t$ -distribution for inference on a sample mean

To proceed with the  $t$ -distribution for inference about a single mean, we first check two conditions.

**Independence of observations.** We verify this condition just as we did before. We collect a simple random sample from less than 10% of the population, or if the data are from an experiment or random process, we check to the best of our abilities that the observations were independent.

**Observations come from a nearly normal distribution.** This second condition is difficult to verify with small data sets. We often (i) take a look at a plot of the data for obvious departures from the normal model, and (ii) consider whether any previous experiences alert us that the data may not be nearly normal.

When examining a sample mean and estimated standard error from a sample of  $n$  independent and nearly normal observations, we use a  $t$ -distribution with  $n - 1$  degrees of freedom ( $df$ ). For example, if the sample size was 19, then we would use the  $t$ -distribution with  $df = 19 - 1 = 18$  degrees of freedom and proceed exactly as we did in Chapter 4, except that *now we use the  $t$ -distribution*.

#### TIP: When to use the $t$ -distribution

Use the  $t$ -distribution for inference of the sample mean when observations are independent and nearly normal. You may relax the nearly normal condition as the sample size increases. For example, the data distribution may be moderately skewed when the sample size is at least 30.

<sup>2</sup>We find the shaded area *above* -1.79 (we leave the picture to you). The small left tail is between 0.025 and 0.05, so the larger upper region must have an area between 0.95 and 0.975.

### 5.1.4 One sample $t$ -confidence intervals

Dolphins are at the top of the oceanic food chain, which causes dangerous substances such as mercury to concentrate in their organs and muscles. This is an important problem for both dolphins and other animals, like humans, who occasionally eat them. For instance, this is particularly relevant in Japan where school meals have included dolphin at times.



Figure 5.6: A Risso's dolphin.

Photo by Mike Baird ([www.bairdphotos.com](http://www.bairdphotos.com)). CC BY 2.0 license.

Here we identify a confidence interval for the average mercury content in dolphin muscle using a sample of 19 Risso's dolphins from the Taiji area in Japan.<sup>3</sup> The data are summarized in Table 5.7. The minimum and maximum observed values can be used to evaluate whether or not there are obvious outliers or skew.

$n$	$\bar{x}$	$s$	minimum	maximum
19	4.4	2.3	1.7	9.2

Table 5.7: Summary of mercury content in the muscle of 19 Risso's dolphins from the Taiji area. Measurements are in  $\mu\text{g}/\text{wet g}$  (micrograms of mercury per wet gram of muscle).

- **Example 5.5** Are the independence and normality conditions satisfied for this data set?

The observations are a simple random sample and consist of less than 10% of the population, therefore independence is reasonable. The summary statistics in Table 5.7 do not suggest any skew or outliers; all observations are within 2.5 standard deviations of the mean. Based on this evidence, the normality assumption seems reasonable.

---

<sup>3</sup>Taiji was featured in the movie *The Cove*, and it is a significant source of dolphin and whale meat in Japan. Thousands of dolphins pass through the Taiji area annually, and we will assume these 19 dolphins represent a simple random sample from those dolphins. Data reference: Endo T and Haraguchi K. 2009. High mercury levels in hair samples from residents of Taiji, a Japanese whaling town. Marine Pollution Bulletin 60(5):743-747.

In the normal model, we used  $z^*$  and the standard error to determine the width of a confidence interval. We revise the confidence interval formula slightly when using the  $t$ -distribution:

$$\bar{x} \pm t_{df}^* SE$$

 $t_{df}^*$ 

Multiplication factor for  $t$  conf. interval

The sample mean and estimated standard error are computed just as before ( $\bar{x} = 4.4$  and  $SE = s/\sqrt{n} = 0.528$ ). The value  $t_{df}^*$  is a cutoff we obtain based on the confidence level and the  $t$ -distribution with  $df$  degrees of freedom. Before determining this cutoff, we will first need the degrees of freedom.

### Degrees of freedom for a single sample

If the sample has  $n$  observations and we are examining a single mean, then we use the  $t$ -distribution with  $df = n - 1$  degrees of freedom.

In our current example, we should use the  $t$ -distribution with  $df = 19 - 1 = 18$  degrees of freedom. Then identifying  $t_{18}^*$  is similar to how we found  $z^*$ .

- For a 95% confidence interval, we want to find the cutoff  $t_{18}^*$  such that 95% of the  $t$ -distribution is between  $-t_{18}^*$  and  $t_{18}^*$ .
- We look in the  $t$ -table on page 143, find the column with area totaling 0.05 in the two tails (third column), and then the row with 18 degrees of freedom:  $t_{18}^* = 2.10$ .

Generally the value of  $t_{df}^*$  is slightly larger than what we would get under the normal model with  $z^*$ .

Finally, we can substitute all our values into the confidence interval equation to create the 95% confidence interval for the average mercury content in muscles from Risso's dolphins that pass through the Taiji area:

$$\bar{x} \pm t_{18}^* SE \rightarrow 4.4 \pm 2.10 \times 0.528 \rightarrow (3.29, 5.51)$$

We are 95% confident the average mercury content of muscles in Risso's dolphins is between 3.29 and 5.51  $\mu\text{g}/\text{wet gram}$ , which is considered extremely high.

### Finding a $t$ -confidence interval for the mean

Based on a sample of  $n$  independent and nearly normal observations, a confidence interval for the population mean is

$$\bar{x} \pm t_{df}^* SE$$

where  $\bar{x}$  is the sample mean,  $t_{df}^*$  corresponds to the confidence level and degrees of freedom, and  $SE$  is the standard error as estimated by the sample.

Ⓐ **Guided Practice 5.6** The FDA's webpage provides some data on mercury content of fish.<sup>4</sup> Based on a sample of 15 croaker white fish (Pacific), a sample mean and standard deviation were computed as 0.287 and 0.069 ppm (parts per million), respectively. The 15 observations ranged from 0.18 to 0.41 ppm. We will assume these observations are independent. Based on the summary statistics of the data, do you have any objections to the normality condition of the individual observations?<sup>5</sup>

Ⓑ **Example 5.7** Estimate the standard error of  $\bar{x} = 0.287$  ppm using the data summaries in Guided Practice 5.6. If we are to use the *t*-distribution to create a 90% confidence interval for the actual mean of the mercury content, identify the degrees of freedom we should use and also find  $t_{df}^*$ .

---

The standard error:  $SE = \frac{0.069}{\sqrt{15}} = 0.0178$ . Degrees of freedom:  $df = n - 1 = 14$ .

Looking in the column where two tails is 0.100 (for a 90% confidence interval) and row  $df = 14$ , we identify  $t_{14}^* = 1.76$ .

Ⓐ **Guided Practice 5.8** Using the results of Guided Practice 5.6 and Example 5.7, compute a 90% confidence interval for the average mercury content of croaker white fish (Pacific).<sup>6</sup>

### 5.1.5 One sample *t*-tests

Is the typical US runner getting faster or slower over time? We consider this question in the context of the Cherry Blossom Race, which is a 10-mile race in Washington, DC each spring.<sup>7</sup>

The average time for all runners who finished the Cherry Blossom Race in 2006 was 93.29 minutes (93 minutes and about 17 seconds). We want to determine using data from 100 participants in the 2012 Cherry Blossom Race whether runners in this race are getting faster or slower, versus the other possibility that there has been no change.

Ⓐ **Guided Practice 5.9** What are appropriate hypotheses for this context?<sup>8</sup>

Ⓑ **Guided Practice 5.10** The data come from a simple random sample from less than 10% of all participants, so the observations are independent. However, should we be worried about skew in the data? See Figure 5.8 for a histogram of the differences.<sup>9</sup>

With independence satisfied and slight skew not a concern for this large of a sample, we can proceed with performing a hypothesis test using the *t*-distribution.

---

<sup>4</sup>[www.fda.gov/food/foodborneillnesscontaminants/metals/ucm115644.htm](http://www.fda.gov/food/foodborneillnesscontaminants/metals/ucm115644.htm)

<sup>5</sup>There are no obvious outliers; all observations are within 2 standard deviations of the mean. If there is skew, it is not evident. There are no red flags for the normal model based on this (limited) information, and we do not have reason to believe the mercury content is not nearly normal in this type of fish.

<sup>6</sup> $\bar{x} \pm t_{14}^* SE \rightarrow 0.287 \pm 1.76 \times 0.0178 \rightarrow (0.256, 0.318)$ . We are 90% confident that the average mercury content of croaker white fish (Pacific) is between 0.256 and 0.318 ppm.

<sup>7</sup>[www.cherryblossom.org](http://www.cherryblossom.org)

<sup>8</sup> $H_0$ : The average 10 mile run time was the same for 2006 and 2012.  $\mu = 93.29$  minutes.  $H_A$ : The average 10 mile run time for 2012 was *different* than that of 2006.  $\mu \neq 93.29$  minutes.

<sup>9</sup>With a sample of 100, we should only be concerned if there is very strong skew. The histogram of the data suggests, at worst, slight skew.

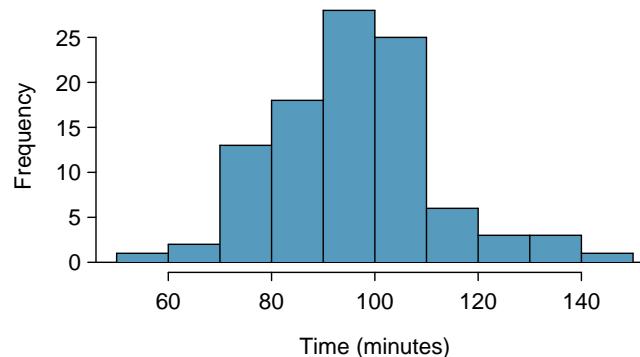


Figure 5.8: A histogram of `time` for the sample Cherry Blossom Race data.

- **Guided Practice 5.11** The sample mean and sample standard deviation of the sample of 100 runners from the 2012 Cherry Blossom Race are 95.61 and 15.78 minutes, respectively. Recall that the sample size is 100. What is the p-value for the test, and what is your conclusion?<sup>10</sup>

#### When using a *t*-distribution, we use a T-score (same as Z-score)

To help us remember to use the *t*-distribution, we use a *T* to represent the test statistic, and we often call this a **T-score**. The Z-score and T-score are computed in the exact same way and are conceptually identical: each represents how many standard errors the observed value is from the null value.

#### 🎥 Calculator videos

Videos covering confidence intervals and hypothesis tests for a single mean using TI and Casio graphing calculators are available at [openintro.org/videos](http://openintro.org/videos).

<sup>10</sup>With the conditions satisfied for the *t*-distribution, we can compute the standard error ( $SE = 15.78/\sqrt{100} = 1.58$  and the *T-score*:  $T = \frac{95.61 - 93.29}{1.58} = 1.47$ ). (There is more on this after the guided practice, but a T-score and Z-score are calculated in the same way.) For  $df = 100 - 1 = 99$ , we would find  $T = 1.47$  to fall between the first and second column, which means the p-value is between 0.10 and 0.20 (use  $df = 90$  and consider two tails since the test is two-sided). The p-value could also have been calculated more precisely with statistical software: 0.1447. Because the p-value is greater than 0.05, we do not reject the null hypothesis. That is, the data do not provide strong evidence that the average run time for the Cherry Blossom Run in 2012 is any different than the 2006 average.

## 5.2 Paired data

Are textbooks actually cheaper online? Here we compare the price of textbooks at the University of California, Los Angeles' (UCLA's) bookstore and prices at Amazon.com. Seventy-three UCLA courses were randomly sampled in Spring 2010, representing less than 10% of all UCLA courses.<sup>11</sup> A portion of the data set is shown in Table 5.9.

	dept	course	ucla	amazon	diff
1	Am Ind	C170	27.67	27.95	-0.28
2	Anthro	9	40.59	31.14	9.45
3	Anthro	135T	31.68	32.00	-0.32
4	Anthro	191HB	16.00	11.52	4.48
:	:	:	:	:	:
72	Wom Std	M144	23.76	18.72	5.04
73	Wom Std	285	27.70	18.22	9.48

Table 5.9: Six cases of the `textbooks` data set.

### 5.2.1 Paired observations

Each textbook has two corresponding prices in the data set: one for the UCLA bookstore and one for Amazon. Therefore, each textbook price from the UCLA bookstore has a natural correspondence with a textbook price from Amazon. When two sets of observations have this special correspondence, they are said to be **paired**.

#### Paired data

Two sets of observations are *paired* if each observation in one set has a special correspondence or connection with exactly one observation in the other data set.

To analyze paired data, it is often useful to look at the difference in outcomes of each pair of observations. In the `textbook` data set, we look at the differences in prices, which is represented as the `diff` variable in the `textbooks` data. Here the differences are taken as

$$\text{UCLA price} - \text{Amazon price}$$

for each book. It is important that we always subtract using a consistent order; here Amazon prices are always subtracted from UCLA prices. A histogram of these differences is shown in Figure 5.10. Using differences between paired observations is a common and useful way to analyze paired data.

- **Guided Practice 5.12** The first difference shown in Table 5.9 is computed as  $27.67 - 27.95 = -0.28$ . Verify the differences are calculated correctly for observations 2 and 3.<sup>12</sup>

<sup>11</sup>When a class had multiple books, only the most expensive text was considered.

<sup>12</sup>Observation 2:  $40.59 - 31.14 = 9.45$ . Observation 3:  $31.68 - 32.00 = -0.32$ .

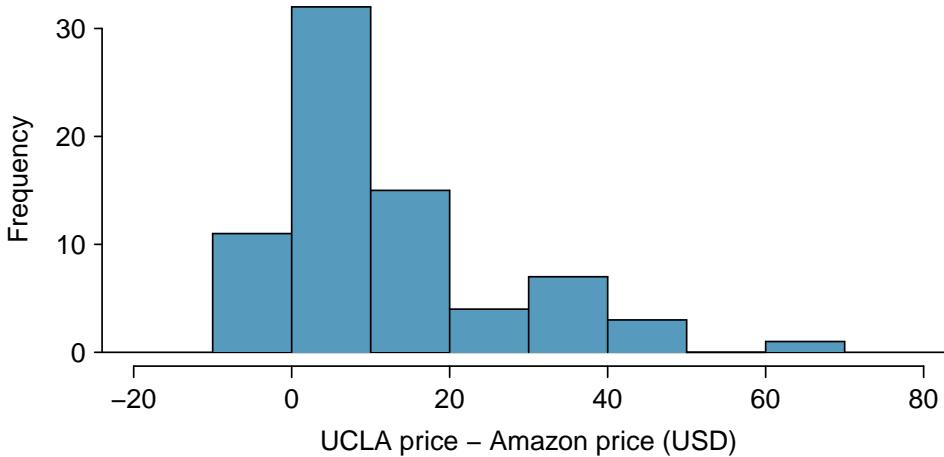


Figure 5.10: Histogram of the difference in price for each book sampled. These data are strongly skewed.

### 5.2.2 Inference for paired data

To analyze a paired data set, we simply analyze the differences. We can use the same *t*-distribution techniques we applied in the last section.

$n_{diff}$	$\bar{x}_{diff}$	$s_{diff}$
73	12.76	14.26

Table 5.11: Summary statistics for the price differences. There were 73 books, so there are 73 differences.

- Example 5.13 Set up and implement a hypothesis test to determine whether, on average, there is a difference between Amazon's price for a book and the UCLA bookstore's price.

We are considering two scenarios: there is no difference or there is some difference in average prices.

$H_0: \mu_{diff} = 0$ . There is no difference in the average textbook price.

$H_A: \mu_{diff} \neq 0$ . There is a difference in average prices.

Can the *t*-distribution be used for this application? The observations are based on a simple random sample from less than 10% of all books sold at the bookstore, so independence is reasonable. While the distribution is strongly skewed, the sample is reasonably large ( $n = 73$ ), so we can proceed. Because the conditions are reasonably satisfied, we can apply the *t*-distribution to this setting.

We compute the standard error associated with  $\bar{x}_{diff}$  using the standard deviation of the differences ( $s_{diff} = 14.26$ ) and the number of differences ( $n_{diff} = 73$ ):

$$SE_{\bar{x}_{diff}} = \frac{s_{diff}}{\sqrt{n_{diff}}} = \frac{14.26}{\sqrt{73}} = 1.67$$

To visualize the p-value, the sampling distribution of  $\bar{x}_{diff}$  is drawn as though  $H_0$  is true, which is shown in Figure 5.12. The p-value is represented by the two (very) small tails.

To find the tail areas, we compute the test statistic, which is the T-score of  $\bar{x}_{diff}$  under the null condition that the actual mean difference is 0:

$$T = \frac{\bar{x}_{diff} - 0}{SE_{\bar{x}_{diff}}} = \frac{12.76 - 0}{1.67} = 7.65$$

The degrees of freedom are  $df = 73 - 1 = 72$ . If we examined Appendix B.2 on page 255, we would see that this value is larger than any in the 70 df row (we round down for  $df$  when using the table), meaning the two-sided p-value is less than 0.01. If we used statistical software, we would find the p-value is less than 1-in-10 billion! Because the p-value is less than 0.05, we reject the null hypothesis. We have found convincing evidence that Amazon was, on average, cheaper than the UCLA bookstore for UCLA course textbooks.

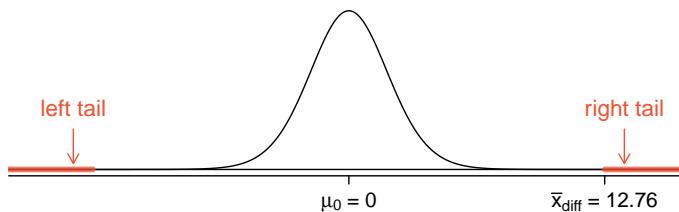


Figure 5.12: Sampling distribution for the mean difference in book prices, if the true average difference is zero.

○ **Guided Practice 5.14** Create a 95% confidence interval for the average price difference between books at the UCLA bookstore and books on Amazon.<sup>13</sup>

## 5.3 Difference of two means

In this section we consider a difference in two population means,  $\mu_1 - \mu_2$ , under the condition that the data are not paired. Just as with a single sample, we identify conditions to ensure we can use the  $t$ -distribution with a point estimate of the difference,  $\bar{x}_1 - \bar{x}_2$ .

We apply these methods in three contexts: determining whether stem cells can improve heart function, exploring the impact of pregnant womens' smoking habits on birth weights of newborns, and exploring whether there is statistically significant evidence that one variation of an exam is harder than another variation. This section is motivated by questions like "Is there convincing evidence that newborns from mothers who smoke have a different average birth weight than newborns from mothers who don't smoke?"

<sup>13</sup>Conditions have already been verified and the standard error computed in Example 5.13. To find the interval, identify  $t_{72}^*$  (use  $df = 70$  in the table,  $t_{70}^* = 1.99$ ) and plug it, the point estimate, and the standard error into the confidence interval formula:

$$\text{point estimate} \pm z^*SE \rightarrow 12.76 \pm 1.99 \times 1.67 \rightarrow (9.44, 16.08)$$

We are 95% confident that Amazon is, on average, between \$9.44 and \$16.08 cheaper than the UCLA bookstore for UCLA course books.

### 5.3.1 Confidence interval for a difference of means

Does treatment using embryonic stem cells (ESCs) help improve heart function following a heart attack? Table 5.13 contains summary statistics for an experiment to test ESCs in sheep that had a heart attack. Each of these sheep was randomly assigned to the ESC or control group, and the change in their hearts' pumping capacity was measured in the study. A positive value corresponds to increased pumping capacity, which generally suggests a stronger recovery. Our goal will be to identify a 95% confidence interval for the effect of ESCs on the change in heart pumping capacity relative to the control group.

A point estimate of the difference in the heart pumping variable can be found using the difference in the sample means:

$$\bar{x}_{esc} - \bar{x}_{control} = 3.50 - (-4.33) = 7.83$$

	$n$	$\bar{x}$	$s$
ESCs	9	3.50	5.17
control	9	-4.33	2.76

Table 5.13: Summary statistics of the embryonic stem cell study.

#### Using the $t$ -distribution for a difference in means

The  $t$ -distribution can be used for inference when working with the standardized difference of two means if (1) each sample meets the conditions for using the  $t$ -distribution and (2) the samples are independent.

- **Example 5.15** Can the  $t$ -distribution be used to make inference using the point estimate,  $\bar{x}_{esc} - \bar{x}_{control} = 7.83$ ?

We check the two required conditions:

1. In this study, the sheep were independent of each other. Additionally, the distributions in Figure 5.14 don't show any clear deviations from normality, where we watch for prominent outliers in particular for such small samples. These findings imply each sample mean could itself be modeled using a  $t$ -distribution.
2. The sheep in each group were also independent of each other.

Because both conditions are met, we can use the  $t$ -distribution to model the difference of the two sample means.

We can quantify the variability in the point estimate,  $\bar{x}_{esc} - \bar{x}_{control}$ , using the following formula for its standard error:

$$SE_{\bar{x}_{esc} - \bar{x}_{control}} = \sqrt{\frac{\sigma_{esc}^2}{n_{esc}} + \frac{\sigma_{control}^2}{n_{control}}}$$

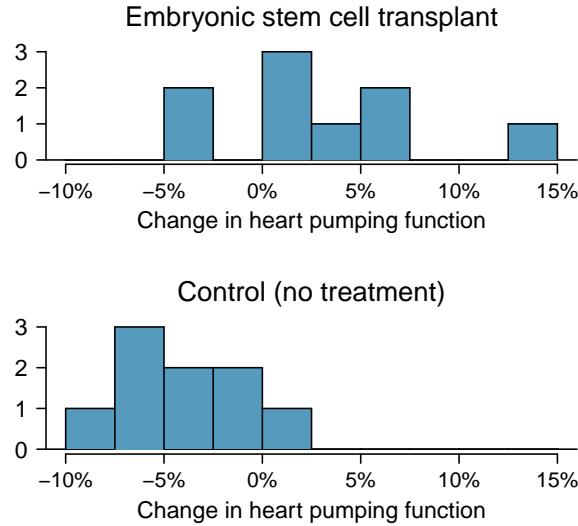


Figure 5.14: Histograms for both the embryonic stem cell group and the control group. Higher values are associated with greater improvement. We don't see any evidence of skew in these data; however, it is worth noting that skew would be difficult to detect with such a small sample.

We usually estimate this standard error using standard deviation estimates based on the samples:

$$\begin{aligned} SE_{\bar{x}_{esc} - \bar{x}_{control}} &= \sqrt{\frac{\sigma_{esc}^2}{n_{esc}} + \frac{\sigma_{control}^2}{n_{control}}} \\ &\approx \sqrt{\frac{s_{esc}^2}{n_{esc}} + \frac{s_{control}^2}{n_{control}}} = \sqrt{\frac{5.17^2}{9} + \frac{2.76^2}{9}} = 1.95 \end{aligned}$$

Because we will use the  $t$ -distribution, we also must identify the appropriate degrees of freedom. This can be done using computer software. An alternative technique is to use the smaller of  $n_1 - 1$  and  $n_2 - 1$ , which is the method we will typically apply in the examples and guided practice.<sup>14</sup>

### Distribution of a difference of sample means

The sample difference of two means,  $\bar{x}_1 - \bar{x}_2$ , can be modeled using the  $t$ -distribution and the standard error

$$SE_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad (5.16)$$

when each sample mean can itself be modeled using a  $t$ -distribution and the samples are independent. To calculate the degrees of freedom, use statistical software or the smaller of  $n_1 - 1$  and  $n_2 - 1$ .

<sup>14</sup>This technique for degrees of freedom is conservative with respect to a Type 1 Error; it is more difficult to reject the null hypothesis using this  $df$  method. In this example, computer software would have provided us a more precise degrees of freedom of  $df = 12.225$ .

- **Example 5.17** Calculate a 95% confidence interval for the effect of ESCs on the change in heart pumping capacity of sheep after they've suffered a heart attack.

We will use the sample difference and the standard error for that point estimate from our earlier calculations:

$$\bar{x}_{esc} - \bar{x}_{control} = 7.83$$

$$SE = \sqrt{\frac{5.17^2}{9} + \frac{2.76^2}{9}} = 1.95$$

Using  $df = 8$ , we can identify the appropriate  $t_{df}^* = t_8^*$  for a 95% confidence interval as 2.31. Finally, we can enter the values into the confidence interval formula:

$$\text{point estimate} \pm t^*SE \rightarrow 7.83 \pm 2.31 \times 1.95 \rightarrow (3.32, 12.34)$$

We are 95% confident that embryonic stem cells improve the heart's pumping function in sheep that have suffered a heart attack by 3.32% to 12.34%.

### 5.3.2 Hypothesis tests based on a difference in means

A data set called `baby_smoke` represents a random sample of 150 cases of mothers and their newborns in North Carolina over a year. Four cases from this data set are represented in Table 5.15. We are particularly interested in two variables: `weight` and `smoke`. The `weight` variable represents the weights of the newborns and the `smoke` variable describes which mothers smoked during pregnancy. We would like to know, is there convincing evidence that newborns from mothers who smoke have a different average birth weight than newborns from mothers who don't smoke? We will use the North Carolina sample to try to answer this question. The smoking group includes 50 cases and the nonsmoking group contains 100 cases, represented in Figure 5.16.

	fAge	mAge	weeks	weight	sexBaby	smoke
1	NA	13	37	5.00	female	nonsmoker
2	NA	14	36	5.88	female	nonsmoker
3	19	15	41	8.13	male	smoker
:	:	:	:	:	:	:
150	45	50	36	9.25	female	nonsmoker

Table 5.15: Four cases from the `baby_smoke` data set. The value “NA”, shown for the first two entries of the first variable, indicates that piece of data is missing.

- **Example 5.18** Set up appropriate hypotheses to evaluate whether there is a relationship between a mother smoking and average birth weight.

The null hypothesis represents the case of no difference between the groups.

$H_0$ : There is no difference in average birth weight for newborns from mothers who did and did not smoke. In statistical notation:  $\mu_n - \mu_s = 0$ , where  $\mu_n$  represents non-smoking mothers and  $\mu_s$  represents mothers who smoked.

$H_A$ : There is some difference in average newborn weights from mothers who did and did not smoke ( $\mu_n - \mu_s \neq 0$ ).

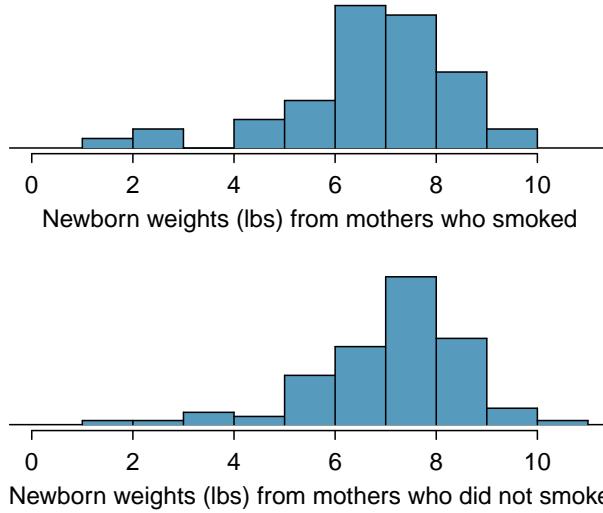


Figure 5.16: The top panel represents birth weights for infants whose mothers smoked. The bottom panel represents the birth weights for infants whose mothers who did not smoke. The distributions exhibit moderate-to-strong skew, respectively.

We check the two conditions necessary to apply the  $t$ -distribution to the difference in sample means. (1) Because the data come from a simple random sample and consist of less than 10% of all such cases, the observations are independent. Additionally, while each distribution is strongly skewed, the sample sizes of 50 and 100 would make it reasonable to model each mean separately using a  $t$ -distribution. The skew is reasonable for these sample sizes of 50 and 100. (2) The independence reasoning applied in (1) also ensures the observations in each sample are independent. Since both conditions are satisfied, the difference in sample means may be modeled using a  $t$ -distribution.

	smoker	nonsmoker
mean	6.78	7.18
st. dev.	1.43	1.60
samp. size	50	100

Table 5.17: Summary statistics for the `baby_smoke` data set.

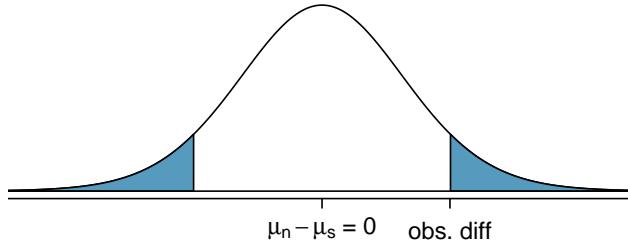
- Ⓐ **Guided Practice 5.19** The summary statistics in Table 5.17 may be useful for this exercise. (a) What is the point estimate of the population difference,  $\mu_n - \mu_s$ ? (b) Compute the standard error of the point estimate from part (a).<sup>15</sup>

<sup>15</sup>(a) The difference in sample means is an appropriate point estimate:  $\bar{x}_n - \bar{x}_s = 0.40$ . (b) The standard error of the estimate can be estimated using Equation (5.16):

$$SE = \sqrt{\frac{\sigma_n^2}{n_n} + \frac{\sigma_s^2}{n_s}} \approx \sqrt{\frac{s_n^2}{n_n} + \frac{s_s^2}{n_s}} = \sqrt{\frac{1.60^2}{100} + \frac{1.43^2}{50}} = 0.26$$

- **Example 5.20** Draw a picture to represent the p-value for the hypothesis test from Example 5.18.

To depict the p-value, we draw the distribution of the point estimate as though  $H_0$  were true and shade areas representing at least as much evidence against  $H_0$  as what was observed. Both tails are shaded because it is a two-sided test.



- **Example 5.21** Compute the p-value of the hypothesis test using the figure in Example 5.20, and evaluate the hypotheses using a significance level of  $\alpha = 0.05$ .

We start by computing the T-score:

$$T = \frac{0.40 - 0}{0.26} = 1.54$$

Next, we compare this value to values in the  $t$ -table in Appendix B.2 on page 255, where we use the smaller of  $n_n - 1 = 99$  and  $n_s - 1 = 49$  as the degrees of freedom:  $df = 49$ . The T-score falls between the first and second columns in the  $df = 49$  row of the  $t$ -table, meaning the two-sided p-value falls between 0.10 and 0.20 (reminder, find tail areas along the top of the table). This p-value is larger than the significance value, 0.05, so we fail to reject the null hypothesis. There is insufficient evidence to say there is a difference in average birth weight of newborns from North Carolina mothers who did smoke during pregnancy and newborns from North Carolina mothers who did not smoke during pregnancy.

- **Guided Practice 5.22** Does the conclusion to Example 5.21 mean that smoking and average birth weight are unrelated?<sup>16</sup>
- **Guided Practice 5.23** If we made a Type 2 Error and there is a difference, what could we have done differently in data collection to be more likely to detect the difference?<sup>17</sup>

<sup>16</sup>Absolutely not. It is possible that there is some difference but we did not detect it. If there is a difference, we made a Type 2 Error. Notice: we also don't have enough information to, if there is an actual difference, confidently say which direction that difference would be in.

<sup>17</sup>We could have collected more data. If the sample sizes are larger, we tend to have a better shot at finding a difference if one exists.

Public service announcement: while we have used this relatively small data set as an example, larger data sets show that women who smoke tend to have smaller newborns. In fact, some in the tobacco industry actually had the audacity to tout that as a *benefit* of smoking:

*It's true. The babies born from women who smoke are smaller, but they're just as healthy as the babies born from women who do not smoke. And some women would prefer having smaller babies.*

- Joseph Cullman, Philip Morris' Chairman of the Board  
on CBS' *Face the Nation*, Jan 3, 1971

Fact check: the babies from women who smoke are not actually as healthy as the babies from women who do not smoke.<sup>18</sup>

### 5.3.3 Case study: two versions of a course exam

An instructor decided to run two slight variations of the same exam. Prior to passing out the exams, she shuffled the exams together to ensure each student received a random version. Summary statistics for how students performed on these two exams are shown in Table 5.18. Anticipating complaints from students who took Version B, she would like to evaluate whether the difference observed in the groups is so large that it provides convincing evidence that Version B was more difficult (on average) than Version A.

Version	$n$	$\bar{x}$	$s$	min	max
A	30	79.4	14	45	100
B	27	74.1	20	32	100

Table 5.18: Summary statistics of scores for each exam version.

Ⓐ **Guided Practice 5.24** Construct a hypotheses to evaluate whether the observed difference in sample means,  $\bar{x}_A - \bar{x}_B = 5.3$ , is due to chance.<sup>19</sup>

Ⓑ **Guided Practice 5.25** To evaluate the hypotheses in Guided Practice 5.24 using the  $t$ -distribution, we must first verify assumptions. (a) Does it seem reasonable that the scores are independent within each group? (b) What about the normality / skew condition for observations in each group? (c) Do you think scores from the two groups would be independent of each other, i.e. the two samples are independent?<sup>20</sup>

After verifying the conditions for each sample and confirming the samples are independent of each other, we are ready to conduct the test using the  $t$ -distribution. In this case,

<sup>18</sup>You can watch an episode of John Oliver on *This Week Tonight* to explore the present day offenses of the tobacco industry. Please be aware that there is some adult language: [youtu.be/6UsHHOCH4q8](https://youtu.be/6UsHHOCH4q8).

<sup>19</sup>Because the teacher did not expect one exam to be more difficult prior to examining the test results, she should use a two-sided hypothesis test.  $H_0$ : the exams are equally difficult, on average.  $\mu_A - \mu_B = 0$ .  $H_A$ : one exam was more difficult than the other, on average.  $\mu_A - \mu_B \neq 0$ .

<sup>20</sup>(a) It is probably reasonable to conclude the scores are independent, provided there was no cheating. (b) The summary statistics suggest the data are roughly symmetric about the mean, and it doesn't seem unreasonable to suggest the data might be normal. Note that since these samples are each nearing 30, moderate skew in the data would be acceptable. (c) It seems reasonable to suppose that the samples are independent since the exams were handed out randomly.

we are estimating the true difference in average test scores using the sample data, so the point estimate is  $\bar{x}_A - \bar{x}_B = 5.3$ . The standard error of the estimate can be calculated as

$$SE = \sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}} = \sqrt{\frac{14^2}{30} + \frac{20^2}{27}} = 4.62$$

Finally, we construct the test statistic:

$$T = \frac{\text{point estimate} - \text{null value}}{SE} = \frac{(79.4 - 74.1) - 0}{4.62} = 1.15$$

If we have a computer handy, we can identify the degrees of freedom as 45.97. Otherwise we use the smaller of  $n_1 - 1$  and  $n_2 - 1$ :  $df = 26$ .

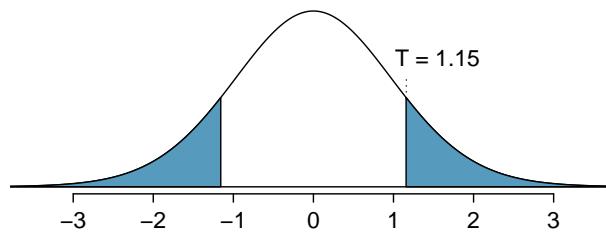


Figure 5.19: The  $t$ -distribution with 26 degrees of freedom. The shaded right tail represents values with  $T \geq 1.15$ . Because it is a two-sided test, we also shade the corresponding lower tail.

- **Example 5.26** Identify the p-value using  $df = 26$  and provide a conclusion in the context of the case study.

We examine row  $df = 26$  in the  $t$ -table. Because this value is smaller than the value in the left column, the p-value is larger than 0.200 (two tails!). Because the p-value is so large, we do not reject the null hypothesis. That is, the data do not convincingly show that one exam version is more difficult than the other, and the teacher should not be convinced that she should add points to the Version B exam scores.

### 5.3.4 Summary for inference using the $t$ -distribution

**Hypothesis tests.** When applying the  $t$ -distribution for a hypothesis test, we proceed as follows:

- Write appropriate hypotheses.
- Verify conditions for using the  $t$ -distribution.
  - One-sample or differences from paired data: the observations (or differences) must be independent and nearly normal. For larger sample sizes, we can relax the nearly normal requirement, e.g. slight skew is okay for sample sizes of 15, moderate skew for sample sizes of 30, and strong skew for sample sizes of 60.
  - For a difference of means when the data are not paired: each sample mean must separately satisfy the one-sample conditions for the  $t$ -distribution, and the data in the groups must also be independent.

- Compute the point estimate of interest, the standard error, and the degrees of freedom. For  $df$ , use  $n - 1$  for one sample, and for two samples use either statistical software or the smaller of  $n_1 - 1$  and  $n_2 - 1$ .
- Compute the T-score and p-value.
- Make a conclusion based on the p-value, and write a conclusion in context and in plain language so anyone can understand the result.

**Confidence intervals.** Similarly, the following is how we generally computed a confidence interval using a  $t$ -distribution:

- Verify conditions for using the  $t$ -distribution. (See above.)
- Compute the point estimate of interest, the standard error, the degrees of freedom, and  $t_{df}^*$ .
- Calculate the confidence interval using the general formula, point estimate  $\pm t_{df}^* SE$ .
- Put the conclusions in context and in plain language so even non-statisticians can understand the results.



### Calculator videos

Videos covering confidence intervals and hypothesis tests for a difference of means using TI and Casio graphing calculators are available at [openintro.org/videos](http://openintro.org/videos).

#### 5.3.5 Examining the standard error formula (special topic)

The formula for the standard error of the difference in two means is similar to the formula for other standard errors. Recall that the standard error of a single mean,  $\bar{x}_1$ , can be approximated by

$$SE_{\bar{x}_1} = \frac{s_1}{\sqrt{n_1}}$$

where  $s_1$  and  $n_1$  represent the sample standard deviation and sample size.

The standard error of the difference of two sample means can be constructed from the standard errors of the separate sample means:

$$SE_{\bar{x}_1 - \bar{x}_2} = \sqrt{SE_{\bar{x}_1}^2 + SE_{\bar{x}_2}^2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad (5.27)$$

This special relationship follows from probability theory.

**Guided Practice 5.28** Prerequisite: Section ?? . We can rewrite Equation (5.27) in a different way:

$$SE_{\bar{x}_1 - \bar{x}_2}^2 = SE_{\bar{x}_1}^2 + SE_{\bar{x}_2}^2$$

Explain where this formula comes from using the ideas of probability theory.<sup>21</sup>

---

<sup>21</sup>The standard error squared represents the variance of the estimate. If  $X$  and  $Y$  are two random variables with variances  $\sigma_x^2$  and  $\sigma_y^2$ , then the variance of  $X - Y$  is  $\sigma_x^2 + \sigma_y^2$ . Likewise, the variance corresponding to  $\bar{x}_1 - \bar{x}_2$  is  $\sigma_{\bar{x}_1}^2 + \sigma_{\bar{x}_2}^2$ . Because  $\sigma_{\bar{x}_1}^2$  and  $\sigma_{\bar{x}_2}^2$  are just another way of writing  $SE_{\bar{x}_1}^2$  and  $SE_{\bar{x}_2}^2$ , the variance associated with  $\bar{x}_1 - \bar{x}_2$  may be written as  $SE_{\bar{x}_1}^2 + SE_{\bar{x}_2}^2$ .

### 5.3.6 Pooled standard deviation estimate (special topic)

Occasionally, two populations will have standard deviations that are so similar that they can be treated as identical. For example, historical data or a well-understood biological mechanism may justify this strong assumption. In such cases, we can make the  $t$ -distribution approach slightly more precise by using a pooled standard deviation.

The **pooled standard deviation** of two groups is a way to use data from both samples to better estimate the standard deviation and standard error. If  $s_1$  and  $s_2$  are the standard deviations of groups 1 and 2 and there are good reasons to believe that the population standard deviations are equal, then we can obtain an improved estimate of the group variances by pooling their data:

$$s_{\text{pooled}}^2 = \frac{s_1^2 \times (n_1 - 1) + s_2^2 \times (n_2 - 1)}{n_1 + n_2 - 2}$$

where  $n_1$  and  $n_2$  are the sample sizes, as before. To use this new statistic, we substitute  $s_{\text{pooled}}^2$  in place of  $s_1^2$  and  $s_2^2$  in the standard error formula, and we use an updated formula for the degrees of freedom:

$$df = n_1 + n_2 - 2$$

The benefits of pooling the standard deviation are realized through obtaining a better estimate of the standard deviation for each group and using a larger degrees of freedom parameter for the  $t$ -distribution. Both of these changes may permit a more accurate model of the sampling distribution of  $\bar{x}_1 - \bar{x}_2$ , if the standard deviations of the two groups are equal.

**Caution: Pool standard deviations only after careful consideration**

A pooled standard deviation is only appropriate when background research indicates the population standard deviations are nearly equal. When the sample size is large and the condition may be adequately checked with data, the benefits of pooling the standard deviations greatly diminishes.

## 5.4 Comparing many means with ANOVA (special topic)



Sometimes we want to compare means across many groups. We might initially think to do pairwise comparisons; for example, if there were three groups, we might be tempted to compare the first mean with the second, then with the third, and then finally compare the second and third means for a total of three comparisons. However, this strategy can be treacherous. If we have many groups and do many comparisons, it is likely that we will eventually find a difference just by chance, even if there is no difference in the populations.

In this section, we will learn a new method called **analysis of variance (ANOVA)** and a new test statistic called  $F$ . ANOVA uses a single hypothesis test to check whether the means across many groups are equal:

$H_0$ : The mean outcome is the same across all groups. In statistical notation,  $\mu_1 = \mu_2 = \dots = \mu_k$  where  $\mu_i$  represents the mean of the outcome for observations in category  $i$ .

$H_A$ : At least one mean is different.

Generally we must check three conditions on the data before performing ANOVA:

- the observations are independent within and across groups,
- the data within each group are nearly normal, and
- the variability across the groups is about equal.

When these three conditions are met, we may perform an ANOVA to determine whether the data provide strong evidence against the null hypothesis that all the  $\mu_i$  are equal.

● **Example 5.29** College departments commonly run multiple lectures of the same introductory course each semester because of high demand. Consider a statistics department that runs three lectures of an introductory statistics course. We might like to determine whether there are statistically significant differences in first exam scores in these three classes ( $A$ ,  $B$ , and  $C$ ). Describe appropriate hypotheses to determine whether there are any differences between the three classes.

---

The hypotheses may be written in the following form:

$H_0$ : The average score is identical in all lectures. Any observed difference is due to chance. Notationally, we write  $\mu_A = \mu_B = \mu_C$ .

$H_A$ : The average score varies by class. We would reject the null hypothesis in favor of the alternative hypothesis if there were larger differences among the class averages than what we might expect from chance alone.

Strong evidence favoring the alternative hypothesis in ANOVA is described by unusually large differences among the group means. We will soon learn that assessing the variability of the group means relative to the variability among individual observations within each group is key to ANOVA's success.

- **Example 5.30** Examine Figure 5.20. Compare groups I, II, and III. Can you visually determine if the differences in the group centers is due to chance or not? Now compare groups IV, V, and VI. Do these differences appear to be due to chance?

Any real difference in the means of groups I, II, and III is difficult to discern, because the data within each group are very volatile relative to any differences in the average outcome. On the other hand, it appears there are differences in the centers of groups IV, V, and VI. For instance, group V appears to have a higher mean than that of the other two groups. Investigating groups IV, V, and VI, we see the differences in the groups' centers are noticeable because those differences are large *relative to the variability in the individual observations within each group*.

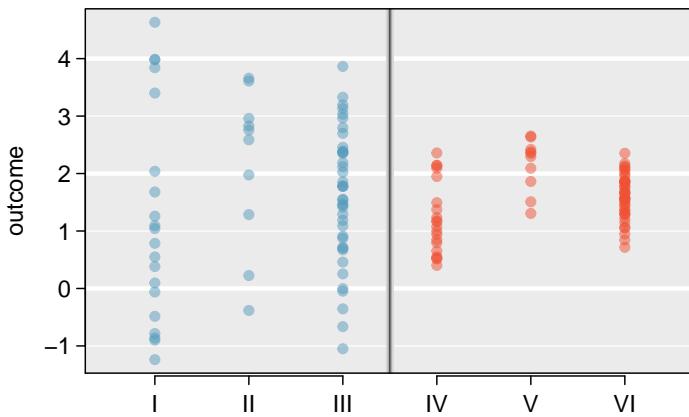


Figure 5.20: Side-by-side dot plot for the outcomes for six groups.

### 5.4.1 Is batting performance related to player position in MLB?

We would like to discern whether there are real differences between the batting performance of baseball players according to their position: outfielder (OF), infielder (IF), designated hitter (DH), and catcher (C). We will use a data set called `bat10`, which includes batting records of 327 Major League Baseball (MLB) players from the 2010 season. Six of the 327 cases represented in `bat10` are shown in Table 5.21, and descriptions for each variable are provided in Table 5.22. The measure we will use for the player batting performance (the outcome variable) is on-base percentage (OBP). The on-base percentage roughly represents the fraction of the time a player successfully gets on base or hits a home run.

- **Guided Practice 5.31** The null hypothesis under consideration is the following:  $\mu_{\text{OF}} = \mu_{\text{IF}} = \mu_{\text{DH}} = \mu_{\text{C}}$ . Write the null and corresponding alternative hypotheses in plain language.<sup>22</sup>

<sup>22</sup> $H_0$ : The average on-base percentage is equal across the four positions.  $H_A$ : The average on-base percentage varies across some (or all) groups.

	name	team	position	AB	H	HR	RBI	AVG	OBP
1	I Suzuki	SEA	OF	680	214	6	43	0.315	0.359
2	D Jeter	NYY	IF	663	179	10	67	0.270	0.340
3	M Young	TEX	IF	656	186	21	91	0.284	0.330
:	:	:	:	:	:	:	:	:	:
325	B Molina	SF	C	202	52	3	17	0.257	0.312
326	J Thole	NYM	C	202	56	3	17	0.277	0.357
327	C Heisey	CIN	OF	201	51	8	21	0.254	0.324

Table 5.21: Six cases from the `bat10` data matrix.

variable	description
<code>name</code>	Player name
<code>team</code>	The abbreviated name of the player's team
<code>position</code>	The player's primary field position (OF, IF, DH, C)
<code>AB</code>	Number of opportunities at bat
<code>H</code>	Number of hits
<code>HR</code>	Number of home runs
<code>RBI</code>	Number of runs batted in
<code>AVG</code>	Batting average, which is equal to $H/AB$
<code>OBP</code>	On-base percentage, which is roughly equal to the fraction of times a player gets on base or hits a home run

Table 5.22: Variables and their descriptions for the `bat10` data set.

● **Example 5.32** The player positions have been divided into four groups: outfield (OF), infield (IF), designated hitter (DH), and catcher (C). What would be an appropriate point estimate of the on-base percentage by outfielders,  $\mu_{OF}$ ?

A good estimate of the on-base percentage by outfielders would be the sample average of OBP for just those players whose position is outfield:  $\bar{x}_{OF} = 0.334$ .

Table 5.23 provides summary statistics for each group. A side-by-side box plot for the on-base percentage is shown in Figure 5.24. Notice that the variability appears to be approximately constant across groups; nearly constant variance across groups is an important assumption that must be satisfied before we consider the ANOVA approach.

	OF	IF	DH	C
Sample size ( $n_i$ )	120	154	14	39
Sample mean ( $\bar{x}_i$ )	0.334	0.332	0.348	0.323
Sample SD ( $s_i$ )	0.029	0.037	0.036	0.045

Table 5.23: Summary statistics of on-base percentage, split by player position.

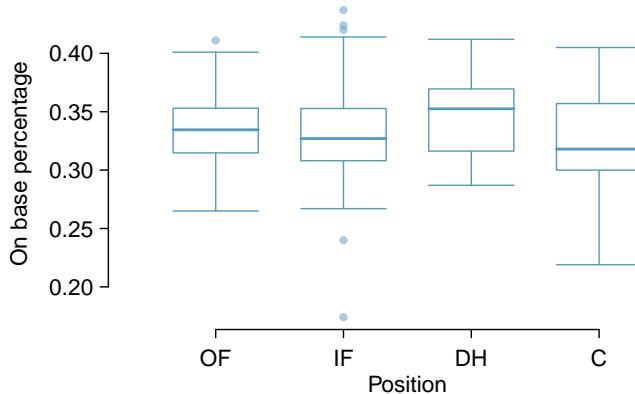


Figure 5.24: Side-by-side box plot of the on-base percentage for 327 players across four groups. There is one prominent outlier visible in the infield group, but with 154 observations in the infield group, this outlier is not a concern.

- ➊ **Example 5.33** The largest difference between the sample means is between the designated hitter and the catcher positions. Consider again the original hypotheses:

$$H_0: \mu_{\text{OF}} = \mu_{\text{IF}} = \mu_{\text{DH}} = \mu_c$$

$H_A$ : The average on-base percentage ( $\mu_i$ ) varies across some (or all) groups.

Why might it be inappropriate to run the test by simply estimating whether the difference of  $\mu_{\text{DH}}$  and  $\mu_c$  is statistically significant at a 0.05 significance level?

---

The primary issue here is that we are inspecting the data before picking the groups that will be compared. It is inappropriate to examine all data by eye (informal testing) and only afterwards decide which parts to formally test. This is called **data snooping** or **data fishing**. Naturally we would pick the groups with the large differences for the formal test, leading to an inflation in the Type 1 Error rate. To understand this better, let's consider a slightly different problem.

Suppose we are to measure the aptitude for students in 20 classes in a large elementary school at the beginning of the year. In this school, all students are randomly assigned to classrooms, so any differences we observe between the classes at the start of the year are completely due to chance. However, with so many groups, we will probably observe a few groups that look rather different from each other. If we select only these classes that look so different, we will probably make the wrong conclusion that the assignment wasn't random. While we might only formally test differences for a few pairs of classes, we informally evaluated the other classes by eye before choosing the most extreme cases for a comparison.

For additional information on the ideas expressed in Example 5.33, we recommend reading about the **prosecutor's fallacy**.<sup>23</sup>

In the next section we will learn how to use the  $F$  statistic and ANOVA to test whether observed differences in sample means could have happened just by chance even if there was no difference in the respective population means.

---

<sup>23</sup>See, for example, [andrewgelman.com/2007/05/18/the\\_prosecutors](http://andrewgelman.com/2007/05/18/the_prosecutors).

### 5.4.2 Analysis of variance (ANOVA) and the F test

The method of analysis of variance in this context focuses on answering one question: is the variability in the sample means so large that it seems unlikely to be from chance alone? This question is different from earlier testing procedures since we will *simultaneously* consider many groups, and evaluate whether their sample means differ more than we would expect from natural variation. We call this variability the **mean square between groups** (*MSG*), and it has an associated degrees of freedom,  $df_G = k - 1$  when there are  $k$  groups. The *MSG* can be thought of as a scaled variance formula for means. If the null hypothesis is true, any variation in the sample means is due to chance and shouldn't be too large. Details of *MSG* calculations are provided in the footnote,<sup>24</sup> however, we typically use software for these computations.

The mean square between the groups is, on its own, quite useless in a hypothesis test. We need a benchmark value for how much variability should be expected among the sample means if the null hypothesis is true. To this end, we compute a pooled variance estimate, often abbreviated as the **mean square error** (*MSE*), which has an associated degrees of freedom value  $df_E = n - k$ . It is helpful to think of *MSE* as a measure of the variability within the groups. Details of the computations of the *MSE* are provided in the footnote<sup>25</sup> for interested readers.

When the null hypothesis is true, any differences among the sample means are only due to chance, and the *MSG* and *MSE* should be about equal. As a test statistic for ANOVA, we examine the fraction of *MSG* and *MSE*:

$$F = \frac{MSG}{MSE} \quad (5.34)$$

The *MSG* represents a measure of the between-group variability, and *MSE* measures the variability within each of the groups.

- Ⓐ **Guided Practice 5.35** For the baseball data,  $MSG = 0.00252$  and  $MSE = 0.00127$ . Identify the degrees of freedom associated with *MSG* and *MSE* and verify the *F* statistic is approximately 1.994.<sup>26</sup>

---

<sup>24</sup>Let  $\bar{x}$  represent the mean of outcomes across all groups. Then the mean square between groups is computed as

$$MSG = \frac{1}{df_G} SSG = \frac{1}{k-1} \sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2$$

where *SSG* is called the **sum of squares between groups** and  $n_i$  is the sample size of group  $i$ .

<sup>25</sup>Let  $\bar{x}$  represent the mean of outcomes across all groups. Then the **sum of squares total** (*SST*) is computed as

$$SST = \sum_{i=1}^n (x_i - \bar{x})^2$$

where the sum is over all observations in the data set. Then we compute the **sum of squared errors** (*SSE*) in one of two equivalent ways:

$$\begin{aligned} SSE &= SST - SSG \\ &= (n_1 - 1)s_1^2 + (n_2 - 1)s_2^2 + \cdots + (n_k - 1)s_k^2 \end{aligned}$$

where  $s_i^2$  is the sample variance (square of the standard deviation) of the residuals in group  $i$ . Then the *MSE* is the standardized form of *SSE*:  $MSE = \frac{1}{df_E} SSE$ .

<sup>26</sup>There are  $k = 4$  groups, so  $df_G = k - 1 = 3$ . There are  $n = n_1 + n_2 + n_3 + n_4 = 327$  total observations, so  $df_E = n - k = 323$ . Then the *F* statistic is computed as the ratio of *MSG* and *MSE*:  $F = \frac{MSG}{MSE} = \frac{0.00252}{0.00127} = 1.984 \approx 1.994$ . ( $F = 1.994$  was computed by using values for *MSG* and *MSE* that were not rounded.)

We can use the  $F$  statistic to evaluate the hypotheses in what is called an  **$F$  test**. A p-value can be computed from the  $F$  statistic using an  $F$  distribution, which has two associated parameters:  $df_1$  and  $df_2$ . For the  $F$  statistic in ANOVA,  $df_1 = df_G$  and  $df_2 = df_E$ . An  $F$  distribution with 3 and 323 degrees of freedom, corresponding to the  $F$  statistic for the baseball hypothesis test, is shown in Figure 5.25.

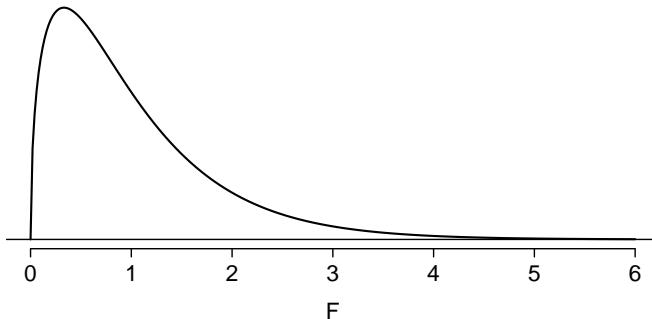


Figure 5.25: An  $F$  distribution with  $df_1 = 3$  and  $df_2 = 323$ .

The larger the observed variability in the sample means ( $MSG$ ) relative to the within-group observations ( $MSE$ ), the larger  $F$  will be and the stronger the evidence against the null hypothesis. Because larger values of  $F$  represent stronger evidence against the null hypothesis, we use the upper tail of the distribution to compute a p-value.

#### The $F$ statistic and the $F$ test

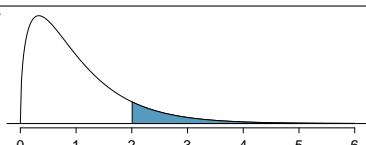
Analysis of variance (ANOVA) is used to test whether the mean outcome differs across 2 or more groups. ANOVA uses a test statistic  $F$ , which represents a standardized ratio of variability in the sample means relative to the variability within the groups. If  $H_0$  is true and the model assumptions are satisfied, the statistic  $F$  follows an  $F$  distribution with parameters  $df_1 = k - 1$  and  $df_2 = n - k$ . The upper tail of the  $F$  distribution is used to represent the p-value.

- Ⓐ **Guided Practice 5.36** The test statistic for the baseball example is  $F = 1.994$ . Shade the area corresponding to the p-value in Figure 5.25.<sup>27</sup>

- Ⓑ **Example 5.37** The p-value corresponding to the shaded area in the solution of Guided Practice 5.36 is equal to about 0.115. Does this provide strong evidence against the null hypothesis?

The p-value is larger than 0.05, indicating the evidence is not strong enough to reject the null hypothesis at a significance level of 0.05. That is, the data do not provide strong evidence that the average on-base percentage varies by player's primary field position.

<sup>27</sup>



### 5.4.3 Reading an ANOVA table from software

The calculations required to perform an ANOVA by hand are tedious and prone to human error. For these reasons, it is common to use statistical software to calculate the  $F$  statistic and p-value.

An ANOVA can be summarized in a table very similar to that of a regression summary, which we will see in Chapters 7 and ???. Table 5.26 shows an ANOVA summary to test whether the mean of on-base percentage varies by player positions in the MLB. Many of these values should look familiar; in particular, the  $F$  test statistic and p-value can be retrieved from the last columns.

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
position	3	0.0076	0.0025	1.9943	0.1147
Residuals	323	0.4080	0.0013		

$$s_{\text{pooled}} = 0.036 \text{ on } df = 323$$

Table 5.26: ANOVA summary for testing whether the average on-base percentage differs across player positions.

### 5.4.4 Graphical diagnostics for an ANOVA analysis

There are three conditions we must check for an ANOVA analysis: all observations must be independent, the data in each group must be nearly normal, and the variance within each group must be approximately equal.

**Independence.** If the data are a simple random sample from less than 10% of the population, this condition is satisfied. For processes and experiments, carefully consider whether the data may be independent (e.g. no pairing). For example, in the MLB data, the data were not sampled. However, there are not obvious reasons why independence would not hold for most or all observations.

**Approximately normal.** As with one- and two-sample testing for means, the normality assumption is especially important when the sample size is quite small. The normal probability plots for each group of the MLB data are shown in Figure 5.27; there is some deviation from normality for infielders, but this isn't a substantial concern since there are about 150 observations in that group and the outliers are not extreme. Sometimes in ANOVA there are so many groups or so few observations per group that checking normality for each group isn't reasonable. See the footnote<sup>28</sup> for guidance on how to handle such instances.

**Constant variance.** The last assumption is that the variance in the groups is about equal from one group to the next. This assumption can be checked by examining a side-by-side box plot of the outcomes across the groups, as in Figure 5.24 on page 164. In this case, the variability is similar in the four groups but not identical. We see in Table 5.23 on page 163 that the standard deviation varies a bit from one group to the next. Whether these differences are from natural variation is unclear, so we should report this uncertainty with the final results.

<sup>28</sup>First calculate the **residuals** of the baseball data, which are calculated by taking the observed values and subtracting the corresponding group means. For example, an outfielder with OBP of 0.405 would have a residual of  $0.405 - \bar{x}_{OF} = 0.071$ . Then to check the normality condition, create a normal probability plot using all the residuals simultaneously.

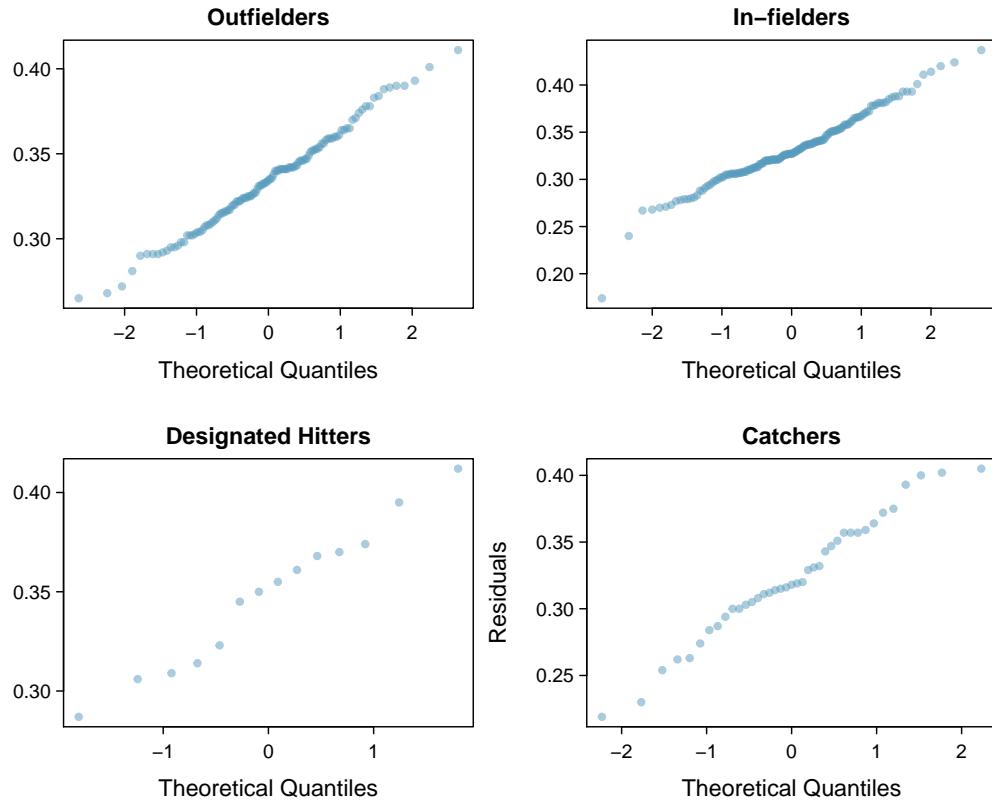


Figure 5.27: Normal probability plot of OBP for each field position.

#### Caution: Diagnostics for an ANOVA analysis

Independence is always important to an ANOVA analysis. The normality condition is very important when the sample sizes for each group are relatively small. The constant variance condition is especially important when the sample sizes differ between groups.

#### 5.4.5 Multiple comparisons and controlling Type 1 Error rate

When we reject the null hypothesis in an ANOVA analysis, we might wonder, which of these groups have different means? To answer this question, we compare the means of each possible pair of groups. For instance, if there are three groups and there is strong evidence that there are some differences in the group means, there are three comparisons to make: group 1 to group 2, group 1 to group 3, and group 2 to group 3. These comparisons can be accomplished using a two-sample  $t$ -test, but we use a modified significance level and a pooled estimate of the standard deviation across groups. Usually this pooled standard deviation can be found in the ANOVA table, e.g. along the bottom of Table 5.26.

Class $i$	A	B	C
$n_i$	58	55	51
$\bar{x}_i$	75.1	72.0	78.9
$s_i$	13.9	13.8	13.1

Table 5.28: Summary statistics for the first midterm scores in three different lectures of the same course.

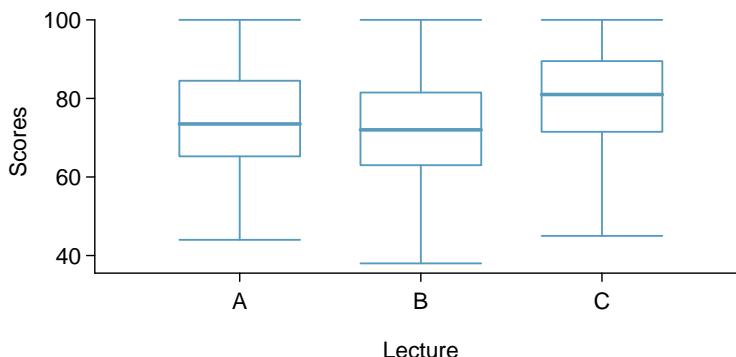


Figure 5.29: Side-by-side box plot for the first midterm scores in three different lectures of the same course.

- **Example 5.38** Example 5.29 on page 161 discussed three statistics lectures, all taught during the same semester. Table 5.28 shows summary statistics for these three courses, and a side-by-side box plot of the data is shown in Figure 5.29. We would like to conduct an ANOVA for these data. Do you see any deviations from the three conditions for ANOVA?

In this case (like many others) it is difficult to check independence in a rigorous way. Instead, the best we can do is use common sense to consider reasons the assumption of independence may not hold. For instance, the independence assumption may not be reasonable if there is a star teaching assistant that only half of the students may access; such a scenario would divide a class into two subgroups. No such situations were evident for these particular data, and we believe that independence is acceptable.

The distributions in the side-by-side box plot appear to be roughly symmetric and show no noticeable outliers.

The box plots show approximately equal variability, which can be verified in Table 5.28, supporting the constant variance assumption.

- **Guided Practice 5.39** An ANOVA was conducted for the midterm data, and summary results are shown in Table 5.30. What should we conclude?<sup>29</sup>

There is strong evidence that the different means in each of the three classes is not simply due to chance. We might wonder, which of the classes are actually different? As

<sup>29</sup>The p-value of the test is 0.0330, less than the default significance level of 0.05. Therefore, we reject the null hypothesis and conclude that the difference in the average midterm scores are not due to chance.

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
lecture	2	1290.11	645.06	3.48	0.0330
Residuals	161	29810.13	185.16		
$s_{pooled} = 13.61$ on $df = 161$					

Table 5.30: ANOVA summary table for the midterm data.

discussed in earlier chapters, a two-sample  $t$ -test could be used to test for differences in each possible pair of groups. However, one pitfall was discussed in Example 5.33 on page 164: when we run so many tests, the Type 1 Error rate increases. This issue is resolved by using a modified significance level.

### Multiple comparisons and the Bonferroni correction for $\alpha$

The scenario of testing many pairs of groups is called **multiple comparisons**. The **Bonferroni correction** suggests that a more stringent significance level is more appropriate for these tests:

$$\alpha^* = \alpha/K$$

where  $K$  is the number of comparisons being considered (formally or informally). If there are  $k$  groups, then usually all possible pairs are compared and  $K = \frac{k(k-1)}{2}$ .

- **Example 5.40** In Guided Practice 5.39, you found strong evidence of differences in the average midterm grades between the three lectures. Complete the three possible pairwise comparisons using the Bonferroni correction and report any differences.

We use a modified significance level of  $\alpha^* = 0.05/3 = 0.0167$ . Additionally, we use the pooled estimate of the standard deviation:  $s_{pooled} = 13.61$  on  $df = 161$ , which is provided in the ANOVA summary table.

Lecture A versus Lecture B: The estimated difference and standard error are, respectively,

$$\bar{x}_A - \bar{x}_B = 75.1 - 72 = 3.1 \quad SE = \sqrt{\frac{13.61^2}{58} + \frac{13.61^2}{55}} = 2.56$$

(See Section 5.3.6 on page 160 for additional details.) This results in a  $T$  score of 1.21 on  $df = 161$  (we use the  $df$  associated with  $s_{pooled}$ ). Statistical software was used to precisely identify the two-sided p-value since the modified significance of 0.0167 is not found in the  $t$ -table. The p-value (0.228) is larger than  $\alpha^* = 0.0167$ , so there is not strong evidence of a difference in the means of lectures A and B.

Lecture A versus Lecture C: The estimated difference and standard error are 3.8 and 2.61, respectively. This results in a  $T$  score of 1.46 on  $df = 161$  and a two-sided p-value of 0.1462. This p-value is larger than  $\alpha^*$ , so there is not strong evidence of a difference in the means of lectures A and C.

Lecture B versus Lecture C: The estimated difference and standard error are 6.9 and 2.65, respectively. This results in a  $T$  score of 2.60 on  $df = 161$  and a two-sided p-value of 0.0102. This p-value is smaller than  $\alpha^*$ . Here we find strong evidence of a difference in the means of lectures B and C.

We might summarize the findings of the analysis from Example 5.40 using the following notation:

$$\mu_A \stackrel{?}{=} \mu_B \quad \mu_A \stackrel{?}{=} \mu_C \quad \mu_B \neq \mu_C$$

The midterm mean in lecture A is not statistically distinguishable from those of lectures B or C. However, there is strong evidence that lectures B and C are different. In the first two pairwise comparisons, we did not have sufficient evidence to reject the null hypothesis. Recall that failing to reject  $H_0$  does not imply  $H_0$  is true.

**Caution: Sometimes an ANOVA will reject the null but no groups will have statistically significant differences**

It is possible to reject the null hypothesis using ANOVA and then to not subsequently identify differences in the pairwise comparisons. However, *this does not invalidate the ANOVA conclusion*. It only means we have not been able to successfully identify which groups differ in their means.

The ANOVA procedure examines the big picture: it considers all groups simultaneously to decipher whether there is evidence that some difference exists. Even if the test indicates that there is strong evidence of differences in group means, identifying with high confidence a specific difference as statistically significant is more difficult.

Consider the following analogy: we observe a Wall Street firm that makes large quantities of money based on predicting mergers. Mergers are generally difficult to predict, and if the prediction success rate is extremely high, that may be considered sufficiently strong evidence to warrant investigation by the Securities and Exchange Commission (SEC). While the SEC may be quite certain that there is insider trading taking place at the firm, the evidence against any single trader may not be very strong. It is only when the SEC considers all the data that they identify the pattern. This is effectively the strategy of ANOVA: stand back and consider all the groups simultaneously.



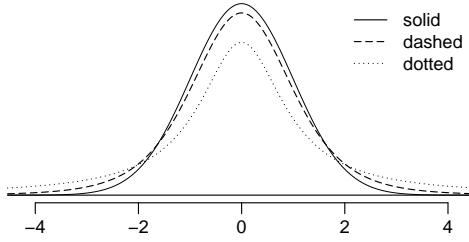
## 5.5 Exercises

### 5.5.1 One-sample means with the t-distribution

**5.1 Identify the critical  $t$ .** An independent random sample is selected from an approximately normal population with unknown standard deviation. Find the degrees of freedom and the critical  $t$ -value ( $t^*$ ) for the given sample size and confidence level.

- (a)  $n = 6$ , CL = 90%
- (b)  $n = 21$ , CL = 98%
- (c)  $n = 29$ , CL = 95%
- (d)  $n = 12$ , CL = 99%

**5.2  $t$ -distribution.** The figure on the right shows three unimodal and symmetric curves: the standard normal (z) distribution, the  $t$ -distribution with 5 degrees of freedom, and the  $t$ -distribution with 1 degree of freedom. Determine which is which, and explain your reasoning.



**5.3 Find the p-value, Part I.** An independent random sample is selected from an approximately normal population with an unknown standard deviation. Find the p-value for the given set of hypotheses and  $T$  test statistic. Also determine if the null hypothesis would be rejected at  $\alpha = 0.05$ .

- (a)  $H_A : \mu > \mu_0$ ,  $n = 11$ ,  $T = 1.91$
- (b)  $H_A : \mu < \mu_0$ ,  $n = 17$ ,  $T = -3.45$
- (c)  $H_A : \mu \neq \mu_0$ ,  $n = 7$ ,  $T = 0.83$
- (d)  $H_A : \mu > \mu_0$ ,  $n = 28$ ,  $T = 2.13$

**5.4 Find the p-value, Part II.** An independent random sample is selected from an approximately normal population with an unknown standard deviation. Find the p-value for the given set of hypotheses and  $T$  test statistic. Also determine if the null hypothesis would be rejected at  $\alpha = 0.01$ .

- (a)  $H_A : \mu > 0.5$ ,  $n = 26$ ,  $T = 2.485$
- (b)  $H_A : \mu < 3$ ,  $n = 18$ ,  $T = 0.5$

**5.5 Working backwards, Part I.** A 95% confidence interval for a population mean,  $\mu$ , is given as (18.985, 21.015). This confidence interval is based on a simple random sample of 36 observations. Calculate the sample mean and standard deviation. Assume that all conditions necessary for inference are satisfied. Use the  $t$ -distribution in any calculations.

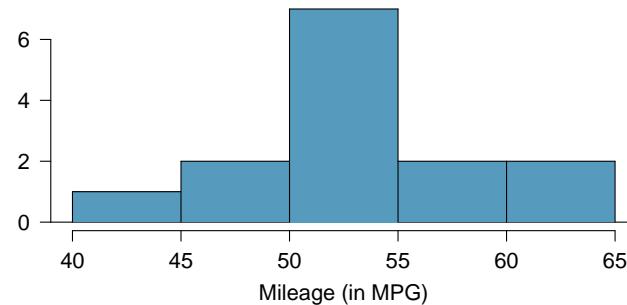
**5.6 Working backwards, Part II.** A 90% confidence interval for a population mean is (65, 77). The population distribution is approximately normal and the population standard deviation is unknown. This confidence interval is based on a simple random sample of 25 observations. Calculate the sample mean, the margin of error, and the sample standard deviation.

**5.7 Sleep habits of New Yorkers.** New York is known as “the city that never sleeps”. A random sample of 25 New Yorkers were asked how much sleep they get per night. Statistical summaries of these data are shown below. Do these data provide strong evidence that New Yorkers sleep less than 8 hours a night on average?

n	$\bar{x}$	s	min	max
25	7.73	0.77	6.17	9.78

- (a) Write the hypotheses in symbols and in words.
- (b) Check conditions, then calculate the test statistic,  $T$ , and the associated degrees of freedom.
- (c) Find and interpret the p-value in this context. Drawing a picture may be helpful.
- (d) What is the conclusion of the hypothesis test?
- (e) If you were to construct a 90% confidence interval that corresponded to this hypothesis test, would you expect 8 hours to be in the interval?

**5.8 Fuel efficiency of Prius.** Fueleconomy.gov, the official US government source for fuel economy information, allows users to share gas mileage information on their vehicles. The histogram below shows the distribution of gas mileage in miles per gallon (MPG) from 14 users who drive a 2012 Toyota Prius. The sample mean is 53.3 MPG and the standard deviation is 5.2 MPG. Note that these data are user estimates and since the source data cannot be verified, the accuracy of these estimates are not guaranteed.<sup>30</sup>



- (a) We would like to use these data to evaluate the average gas mileage of all 2012 Prius drivers. Do you think this is reasonable? Why or why not?
- (b) The EPA claims that a 2012 Prius gets 50 MPG (city and highway mileage combined). Do these data provide strong evidence against this estimate for drivers who participate on fueleconomy.gov? Note any assumptions you must make as you proceed with the test.
- (c) Calculate a 95% confidence interval for the average gas mileage of a 2012 Prius by drivers who participate on fueleconomy.gov.

**5.9 Find the mean.** You are given the following hypotheses:

$$H_0 : \mu = 60$$

$$H_A : \mu < 60$$

We know that the sample standard deviation is 8 and the sample size is 20. For what sample mean would the p-value be equal to 0.05? Assume that all conditions necessary for inference are satisfied.

**5.10  $t^*$  vs.  $z^*$ .** For a given confidence level,  $t_{df}^*$  is larger than  $z^*$ . Explain how  $t_{df}^*$  being slightly larger than  $z^*$  affects the width of the confidence interval.

<sup>30</sup>Fueleconomy.gov, Shared MPG Estimates: Toyota Prius 2012.



**5.11 Play the piano.** Georgianna claims that in a small city renowned for its music school, the average child takes at least 5 years of piano lessons. We have a random sample of 20 children from the city, with a mean of 4.6 years of piano lessons and a standard deviation of 2.2 years.

- Evaluate Georgianna's claim using a hypothesis test.
- Construct a 95% confidence interval for the number of years students in this city take piano lessons, and interpret it in context of the data.
- Do your results from the hypothesis test and the confidence interval agree? Explain your reasoning.

**5.12 Auto exhaust and lead exposure.** Researchers interested in lead exposure due to car exhaust sampled the blood of 52 police officers subjected to constant inhalation of automobile exhaust fumes while working traffic enforcement in a primarily urban environment. The blood samples of these officers had an average lead concentration of  $124.32 \mu\text{g/l}$  and a SD of  $37.74 \mu\text{g/l}$ ; a previous study of individuals from a nearby suburb, with no history of exposure, found an average blood level concentration of  $35 \mu\text{g/l}$ .<sup>31</sup>

- Write down the hypotheses that would be appropriate for testing if the police officers appear to have been exposed to a higher concentration of lead.
- Explicitly state and check all conditions necessary for inference on these data.
- Test the hypothesis that the downtown police officers have a higher lead exposure than the group in the previous study. Interpret your results in context.
- Based on your preceding result, without performing a calculation, would a 99% confidence interval for the average blood concentration level of police officers contain  $35 \mu\text{g/l}$ ?



**5.13 Car insurance savings.** A market researcher wants to evaluate car insurance savings at a competing company. Based on past studies he is assuming that the standard deviation of savings is \$100. He wants to collect data such that he can get a margin of error of no more than \$10 at a 95% confidence level. How large of a sample should he collect?

**5.14 SAT scores.** SAT scores of students at an Ivy League college are distributed with a standard deviation of 250 points. Two statistics students, Raina and Luke, want to estimate the average SAT score of students at this college as part of a class project. They want their margin of error to be no more than 25 points.

- Raina wants to use a 90% confidence interval. How large a sample should she collect?
- Luke wants to use a 99% confidence interval. Without calculating the actual sample size, determine whether his sample should be larger or smaller than Raina's, and explain your reasoning.
- Calculate the minimum required sample size for Luke.

## 5.5.2 Paired data

**5.15 Air quality.** Air quality measurements were collected in a random sample of 25 country capitals in 2013, and then again in the same cities in 2014. We would like to use these data to compare average air quality between the two years.

- Should we use a one-sided or a two-sided test? Explain your reasoning.
- Should we use a paired or non-paired test? Explain your reasoning.
- Should we use a  $t$ -test or a  $z$ -test? Explain your reasoning.

---

<sup>31</sup>WI Mortada et al. "Study of lead exposure from automobile exhaust as a risk for nephrotoxicity among traffic policemen." In: *American journal of nephrology* 21.4 (2000), pp. 274–279.

**5.16 True / False: paired.** Determine if the following statements are true or false. If false, explain.

- (a) In a paired analysis we first take the difference of each pair of observations, and then we do inference on these differences.
- (b) Two data sets of different sizes cannot be analyzed as paired data.
- (c) Consider two data sets that form paired data. Each observation in one data set has a natural correspondence with exactly one observation from the other data set.
- (d) Consider two data sets that form paired data. In the analysis, each observation in one data set is subtracted from the average of the other data set's observations.

**5.17 Paired or not, Part I?** In each of the following scenarios, determine if the data are paired.

- (a) Compare pre- (beginning of semester) and post-test (end of semester) scores of students.
- (b) Assess gender-related salary gap by comparing salaries of randomly sampled men and women.
- (c) Compare artery thicknesses at the beginning of a study and after 2 years of taking Vitamin E for the same group of patients.
- (d) Assess effectiveness of a diet regimen by comparing the before and after weights of subjects.

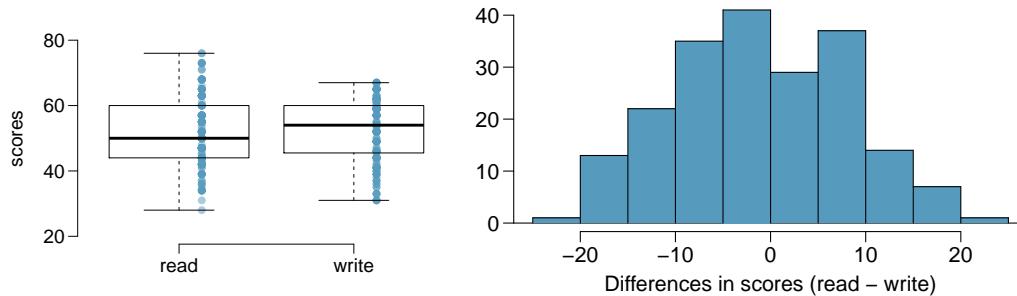
**5.18 Paired or not, Part II?** In each of the following scenarios, determine if the data are paired.

- (a) We would like to know if Intel's stock and Southwest Airlines' stock have similar rates of return. To find out, we take a random sample of 50 days, and record Intel's and Southwest's stock on those same days.
- (b) We randomly sample 50 items from Target stores and note the price for each. Then we visit Walmart and collect the price for each of those same 50 items.
- (c) A school board would like to determine whether there is a difference in average SAT scores for students at one high school versus another high school in the district. To check, they take a simple random sample of 100 students from each high school.

**5.19 Global warming, Part I.** Is there strong evidence of global warming? Let's consider a small scale example, comparing how temperatures have changed in the US from 1968 to 2008. The daily high temperature reading on January 1 was collected in 1968 and 2008 for 51 randomly selected locations in the continental US. Then the difference between the two readings (temperature in 2008 - temperature in 1968) was calculated for each of the 51 different locations. The average of these 51 values was 1.1 degrees with a standard deviation of 4.9 degrees. We are interested in determining whether these data provide strong evidence of temperature warming in the continental US.

- (a) Is there a relationship between the observations collected in 1968 and 2008? Or are the observations in the two groups independent? Explain.
- (b) Write hypotheses for this research in symbols and in words.
- (c) Check the conditions required to complete this test.
- (d) Calculate the test statistic and find the p-value.
- (e) What do you conclude? Interpret your conclusion in context.
- (f) What type of error might we have made? Explain in context what the error means.
- (g) Based on the results of this hypothesis test, would you expect a confidence interval for the average difference between the temperature measurements from 1968 and 2008 to include 0? Explain your reasoning.

**5.20 High School and Beyond, Part I.** The National Center of Education Statistics conducted a survey of high school seniors, collecting test data on reading, writing, and several other subjects. Here we examine a simple random sample of 200 students from this survey. Side-by-side box plots of reading and writing scores as well as a histogram of the differences in scores are shown below.



- (a) Is there a clear difference in the average reading and writing scores?
- (b) Are the reading and writing scores of each student independent of each other?
- (c) Create hypotheses appropriate for the following research question: is there an evident difference in the average scores of students in the reading and writing exam?
- (d) Check the conditions required to complete this test.
- (e) The average observed difference in scores is  $\bar{x}_{\text{read-write}} = -0.545$ , and the standard deviation of the differences is 8.887 points. Do these data provide convincing evidence of a difference between the average scores on the two exams?
- (f) What type of error might we have made? Explain what the error means in the context of the application.
- (g) Based on the results of this hypothesis test, would you expect a confidence interval for the average difference between the reading and writing scores to include 0? Explain your reasoning.

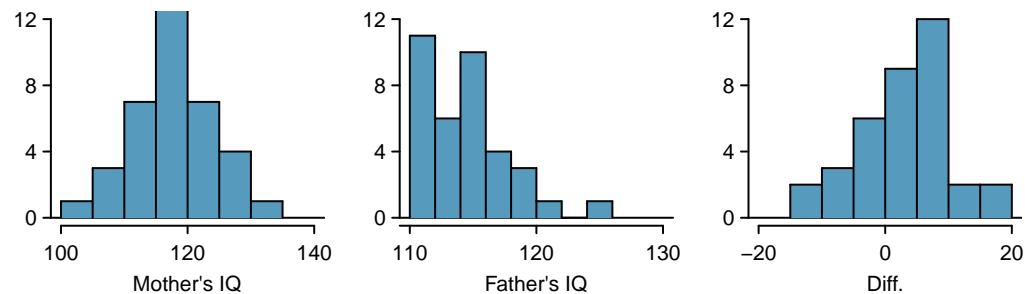
**5.21 Global warming, Part II.** We considered the differences between the temperature readings in January 1 of 1968 and 2008 at 51 locations in the continental US in Exercise 5.19. The mean and standard deviation of the reported differences are 1.1 degrees and 4.9 degrees.

- (a) Calculate a 90% confidence interval for the average difference between the temperature measurements between 1968 and 2008.
- (b) Interpret this interval in context.
- (c) Does the confidence interval provide convincing evidence that the temperature was higher in 2008 than in 1968 in the continental US? Explain.

**5.22 High school and beyond, Part II.** We considered the differences between the reading and writing scores of a random sample of 200 students who took the High School and Beyond Survey in Exercise 5.20. The mean and standard deviation of the differences are  $\bar{x}_{\text{read-write}} = -0.545$  and 8.887 points.

- (a) Calculate a 95% confidence interval for the average difference between the reading and writing scores of all students.
- (b) Interpret this interval in context.
- (c) Does the confidence interval provide convincing evidence that there is a real difference in the average scores? Explain.

**5.23 Gifted children.** Researchers collected a simple random sample of 36 children who had been identified as gifted in a large city. The following histograms show the distributions of the IQ scores of mothers and fathers of these children. Also provided are some sample statistics.<sup>32</sup>



	Mother	Father	Diff.
Mean	118.2	114.8	3.4
SD	6.5	3.5	7.5
n	36	36	36

- (a) Are the IQs of mothers and the IQs of fathers in this data set related? Explain.
- (b) Conduct a hypothesis test to evaluate if the scores are equal on average. Make sure to clearly state your hypotheses, check the relevant conditions, and state your conclusion in the context of the data.

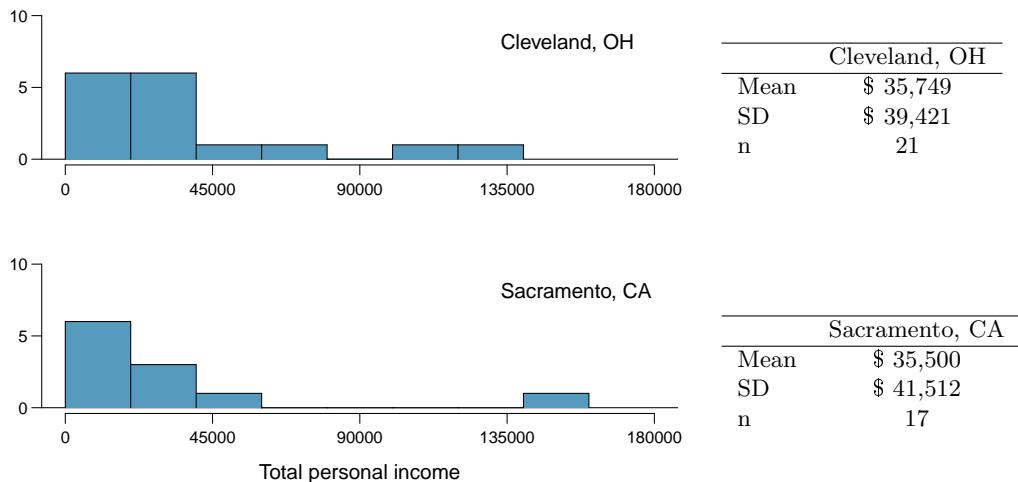
**5.24 Sample size and pairing.** Determine if the following statement is true or false, and if false, explain your reasoning: If comparing means of two groups with equal sample sizes, always use a paired test.

---

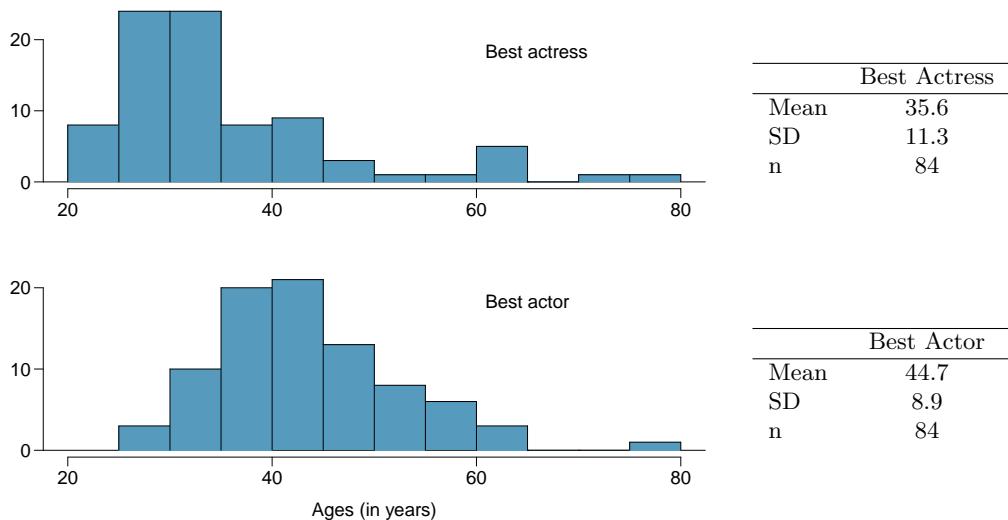
<sup>32</sup>F.A. Graybill and H.K. Iyer. *Regression Analysis: Concepts and Applications*. Duxbury Press, 1994, pp. 511–516.

### 5.5.3 Difference of two means

**5.25 Cleveland vs. Sacramento.** Average income varies from one region of the country to another, and it often reflects both lifestyles and regional living expenses. Suppose a new graduate is considering a job in two locations, Cleveland, OH and Sacramento, CA, and he wants to see whether the average income in one of these cities is higher than the other. He would like to conduct a hypothesis test based on two small samples from the 2000 Census, but he first must consider whether the conditions are met to implement the test. Below are histograms for each city. Should he move forward with the hypothesis test? Explain your reasoning.

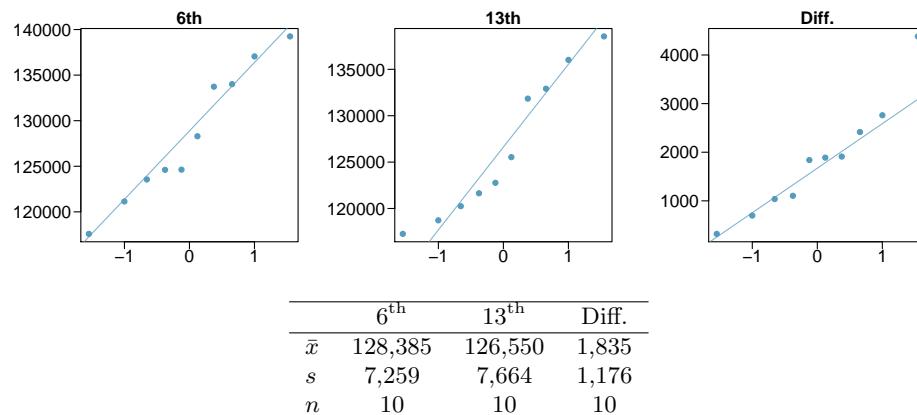


**5.26 Oscar winners.** The first Oscar awards for best actor and best actress were given out in 1929. The histograms below show the age distribution for all of the best actor and best actress winners from 1929 to 2012. Summary statistics for these distributions are also provided. Is a hypothesis test appropriate for evaluating whether the difference in the average ages of best actors and actresses might be due to chance? Explain your reasoning.<sup>33</sup>



<sup>33</sup>Oscar winners from 1929 – 2012, data up to 2009 from the Journal of Statistics Education data archive and more current data from wikipedia.org.

**5.27 Friday the 13<sup>th</sup>, Part I.** In the early 1990's, researchers in the UK collected data on traffic flow, number of shoppers, and traffic accident related emergency room admissions on Friday the 13<sup>th</sup> and the previous Friday, Friday the 6<sup>th</sup>. The histograms below show the distribution of number of cars passing by a specific intersection on Friday the 6<sup>th</sup> and Friday the 13<sup>th</sup> for many such date pairs. Also given are some sample statistics, where the difference is the number of cars on the 6<sup>th</sup> minus the number of cars on the 13<sup>th</sup>.<sup>34</sup>

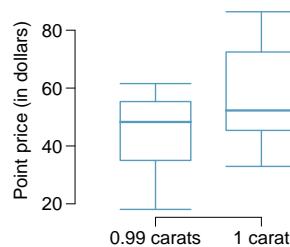


- (a) Are there any underlying structures in these data that should be considered in an analysis? Explain.
- (b) What are the hypotheses for evaluating whether the number of people out on Friday the 6<sup>th</sup> is different than the number out on Friday the 13<sup>th</sup>?
- (c) Check conditions to carry out the hypothesis test from part (b).
- (d) Calculate the test statistic and the p-value.
- (e) What is the conclusion of the hypothesis test?
- (f) Interpret the p-value in this context.
- (g) What type of error might have been made in the conclusion of your test? Explain.

**5.28 Diamonds, Part I.** Prices of diamonds are determined by what is known as the 4 Cs: cut, clarity, color, and carat weight. The prices of diamonds go up as the carat weight increases, but the increase is not smooth. For example, the difference between the size of a 0.99 carat diamond and a 1 carat diamond is undetectable to the naked human eye, but the price of a 1 carat diamond tends to be much higher than the price of a 0.99 diamond. In this question we use two random samples of diamonds, 0.99 carats and 1 carat, each sample of size 23, and compare the average prices of the diamonds. In order to be able to compare equivalent units, we first divide the price for each diamond by 100 times its weight in carats. That is, for a 0.99 carat diamond, we divide the price by 99. For a 1 carat diamond, we divide the price by 100. The distributions and some sample statistics are shown below.<sup>35</sup>

Conduct a hypothesis test to evaluate if there is a difference between the average standardized prices of 0.99 and 1 carat diamonds. Make sure to state your hypotheses clearly, check relevant conditions, and interpret your results in context of the data.

	0.99 carats	1 carat
Mean	\$ 44.51	\$ 56.81
SD	\$ 13.32	\$ 16.13
n	23	23

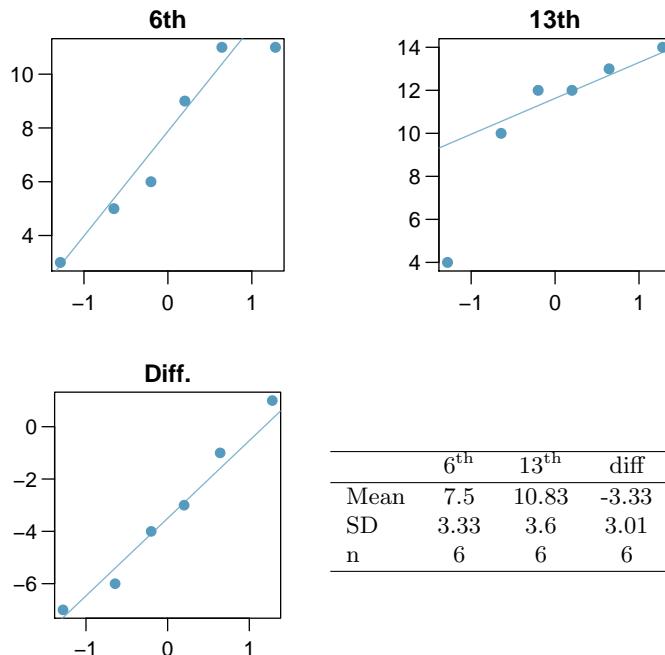


<sup>34</sup>T.J. Scanlon et al. "Is Friday the 13th Bad For Your Health?" In: *BMJ* 307 (1993), pp. 1584–1586.

<sup>35</sup>H. Wickham. *ggplot2: elegant graphics for data analysis*. Springer New York, 2009.



**5.29 Friday the 13<sup>th</sup>, Part II.** The Friday the 13<sup>th</sup> study reported in Exercise 5.27 also provides data on traffic accident related emergency room admissions. The distributions of these counts from Friday the 6<sup>th</sup> and Friday the 13<sup>th</sup> are shown below for six such paired dates along with summary statistics. You may assume that conditions for inference are met.



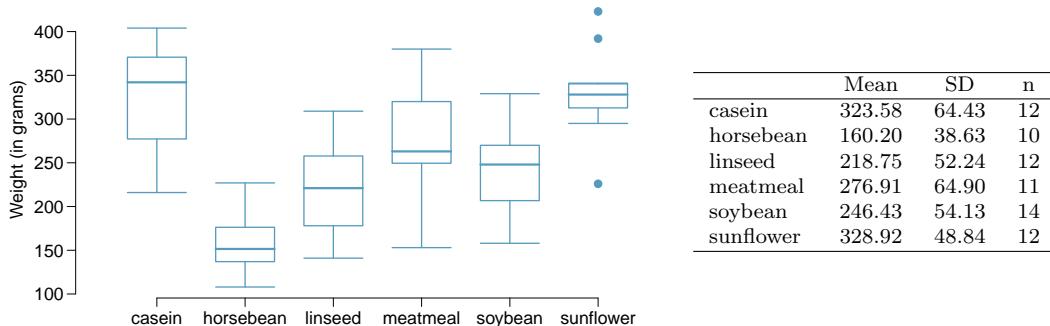
- (a) Conduct a hypothesis test to evaluate if there is a difference between the average numbers of traffic accident related emergency room admissions between Friday the 6<sup>th</sup> and Friday the 13<sup>th</sup>.
- (b) Calculate a 95% confidence interval for the difference between the average numbers of traffic accident related emergency room admissions between Friday the 6<sup>th</sup> and Friday the 13<sup>th</sup>.
- (c) The conclusion of the original study states, “Friday 13th is unlucky for some. The risk of hospital admission as a result of a transport accident may be increased by as much as 52%. Staying at home is recommended.” Do you agree with this statement? Explain your reasoning.

**5.30 Diamonds, Part II.** In Exercise 5.28, we discussed diamond prices (standardized by weight) for diamonds with weights 0.99 carats and 1 carat. See the table for summary statistics, and then construct a 95% confidence interval for the average difference between the standardized prices of 0.99 and 1 carat diamonds. You may assume the conditions for inference are met.

	0.99 carats	1 carat
Mean	\$ 44.51	\$ 56.81
SD	\$ 13.32	\$ 16.13
n	23	23



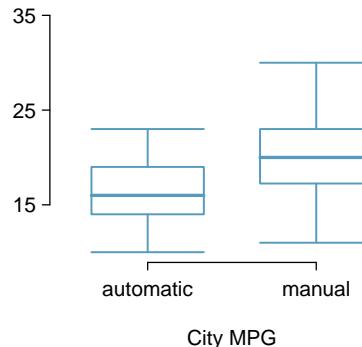
**5.31 Chicken diet and weight, Part I.** Chicken farming is a multi-billion dollar industry, and any methods that increase the growth rate of young chicks can reduce consumer costs while increasing company profits, possibly by millions of dollars. An experiment was conducted to measure and compare the effectiveness of various feed supplements on the growth rate of chickens. Newly hatched chicks were randomly allocated into six groups, and each group was given a different feed supplement. Below are some summary statistics from this data set along with box plots showing the distribution of weights by feed type.<sup>36</sup>



- Describe the distributions of weights of chickens that were fed linseed and horsebean.
- Do these data provide strong evidence that the average weights of chickens that were fed linseed and horsebean are different? Use a 5% significance level.
- What type of error might we have committed? Explain.
- Would your conclusion change if we used  $\alpha = 0.01$ ?

**5.32 Fuel efficiency of manual and automatic cars, Part I.** Each year the US Environmental Protection Agency (EPA) releases fuel economy data on cars manufactured in that year. Below are summary statistics on fuel efficiency (in miles/gallon) from random samples of cars with manual and automatic transmissions manufactured in 2012. Do these data provide strong evidence of a difference between the average fuel efficiency of cars with manual and automatic transmissions in terms of their average city mileage? Assume that conditions for inference are satisfied.<sup>37</sup>

City MPG		
	Automatic	Manual
Mean	16.12	19.85
SD	3.58	4.51
n	26	26



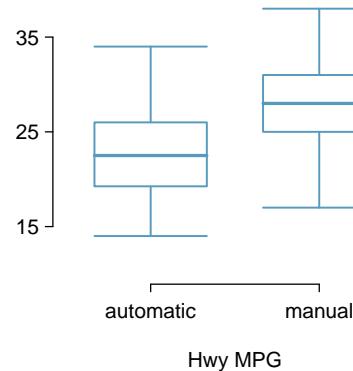
**5.33 Chicken diet and weight, Part II.** Casein is a common weight gain supplement for humans. Does it have an effect on chickens? Using data provided in Exercise 5.31, test the hypothesis that the average weight of chickens that were fed casein is different than the average weight of chickens that were fed soybean. If your hypothesis test yields a statistically significant result, discuss whether or not the higher average weight of chickens can be attributed to the casein diet. Assume that conditions for inference are satisfied.

<sup>36</sup>Chicken Weights by Feed Type, from the `datasets` package in R..

<sup>37</sup>U.S. Department of Energy, Fuel Economy Data, 2012 Datafile.

**5.34 Fuel efficiency of manual and automatic cars, Part II.** The table provides summary statistics on highway fuel economy of cars manufactured in 2012 (from Exercise 5.32). Use these statistics to calculate a 98% confidence interval for the difference between average highway mileage of manual and automatic cars, and interpret this interval in the context of the data.<sup>38</sup>

Hwy MPG		
	Automatic	Manual
Mean	22.92	27.88
SD	5.29	5.01
n	26	26



**5.35 Gaming and distracted eating, Part I.** A group of researchers are interested in the possible effects of distracting stimuli during eating, such as an increase or decrease in the amount of food consumption. To test this hypothesis, they monitored food intake for a group of 44 patients who were randomized into two equal groups. The treatment group ate lunch while playing solitaire, and the control group ate lunch without any added distractions. Patients in the treatment group ate 52.1 grams of biscuits, with a standard deviation of 45.1 grams, and patients in the control group ate 27.1 grams of biscuits, with a standard deviation of 26.4 grams. Do these data provide convincing evidence that the average food intake (measured in amount of biscuits consumed) is different for the patients in the treatment group? Assume that conditions for inference are satisfied.<sup>39</sup>

**5.36 Gaming and distracted eating, Part II.** The researchers from Exercise 5.35 also investigated the effects of being distracted by a game on how much people eat. The 22 patients in the treatment group who ate their lunch while playing solitaire were asked to do a serial-order recall of the food lunch items they ate. The average number of items recalled by the patients in this group was 4.9, with a standard deviation of 1.8. The average number of items recalled by the patients in the control group (no distraction) was 6.1, with a standard deviation of 1.8. Do these data provide strong evidence that the average number of food items recalled by the patients in the treatment and control groups are different?

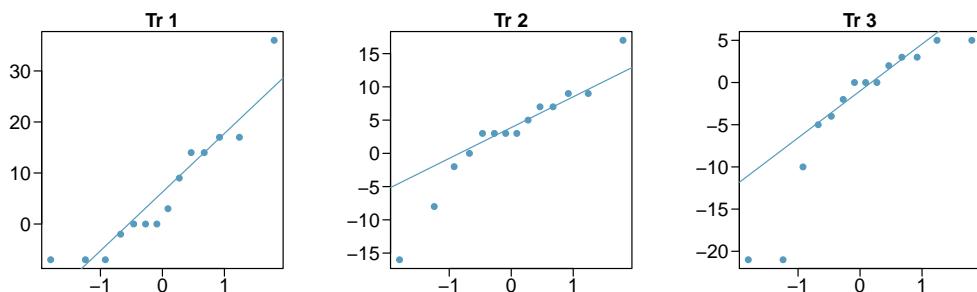
<sup>38</sup>U.S. Department of Energy, Fuel Economy Data, 2012 Datafile.

<sup>39</sup>R.E. Oldham-Cooper et al. “Playing a computer game during lunch affects fullness, memory for lunch, and later snack intake”. In: *The American Journal of Clinical Nutrition* 93.2 (2011), p. 308.

**5.37 Prison isolation experiment, Part I.** Subjects from Central Prison in Raleigh, NC, volunteered for an experiment involving an “isolation” experience. The goal of the experiment was to find a treatment that reduces subjects’ psychopathic deviant T scores. This score measures a person’s need for control or their rebellion against control, and it is part of a commonly used mental health test called the Minnesota Multiphasic Personality Inventory (MMPI) test. The experiment had three treatment groups:

- (1) Four hours of sensory restriction plus a 15 minute “therapeutic” tape advising that professional help is available.
- (2) Four hours of sensory restriction plus a 15 minute “emotionally neutral” tape on training hunting dogs.
- (3) Four hours of sensory restriction but no taped message.

Forty-two subjects were randomly assigned to these treatment groups, and an MMPI test was administered before and after the treatment. Distributions of the differences between pre and post treatment scores (pre - post) are shown below, along with some sample statistics. Use this information to independently test the effectiveness of each treatment. Make sure to clearly state your hypotheses, check conditions, and interpret results in the context of the data.<sup>40</sup>



**5.38 True / False: comparing means.** Determine if the following statements are true or false, and explain your reasoning for statements you identify as false.

- (a) When comparing means of two samples where  $n_1 = 20$  and  $n_2 = 40$ , we can use the normal model for the difference in means since  $n_2 \geq 30$ .
- (b) As the degrees of freedom increases, the  $t$ -distribution approaches normality.
- (c) We use a pooled standard error for calculating the standard error of the difference between means when sample sizes of groups are equal to each other.

---

<sup>40</sup>Prison isolation experiment.

### 5.5.4 Comparing many means with ANOVA

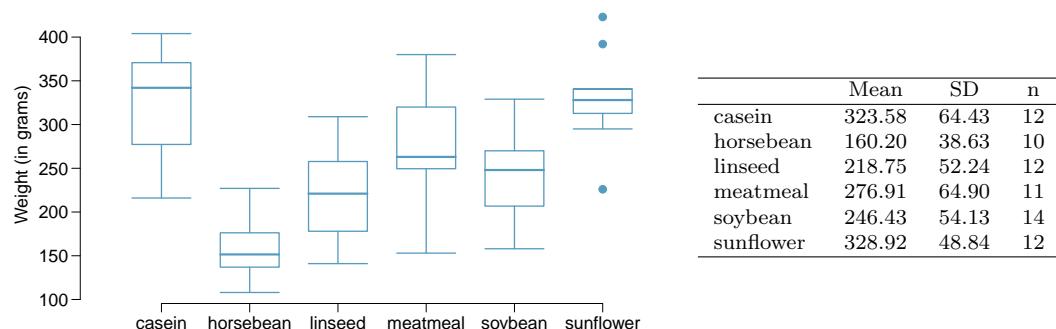
**5.39 Fill in the blank.** When doing an ANOVA, you observe large differences in means between groups. Within the ANOVA framework, this would most likely be interpreted as evidence strongly favoring the \_\_\_\_\_ hypothesis.

**5.40 Which test?** We would like to test if students who are in the social sciences, natural sciences, arts and humanities, and other fields spend the same amount of time studying for this course. What type of test should we use? Explain your reasoning.

**5.41 Chicken diet and weight, Part III.** In Exercises 5.31 and 5.33 we compared the effects of two types of feed at a time. A better analysis would first consider all feed types at once: casein, horsebean, linseed, meat meal, soybean, and sunflower. The ANOVA output below can be used to test for differences between the average weights of chicks on different diets.

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
feed	5	231,129.16	46,225.83	15.36	0.0000
Residuals	65	195,556.02	3,008.55		

Conduct a hypothesis test to determine if these data provide convincing evidence that the average weight of chicks varies across some (or all) groups. Make sure to check relevant conditions. Figures and summary statistics are shown below.



**5.42 Teaching descriptive statistics.** A study compared five different methods for teaching descriptive statistics. The five methods were traditional lecture and discussion, programmed textbook instruction, programmed text with lectures, computer instruction, and computer instruction with lectures. 45 students were randomly assigned, 9 to each method. After completing the course, students took a 1-hour exam.

- (a) What are the hypotheses for evaluating if the average test scores are different for the different teaching methods?
- (b) What are the degrees of freedom associated with the  $F$ -test for evaluating these hypotheses?
- (c) Suppose the p-value for this test is 0.0168. What is the conclusion?

**5.43 Coffee, depression, and physical activity.** Caffeine is the world's most widely used stimulant, with approximately 80% consumed in the form of coffee. Participants in a study investigating the relationship between coffee consumption and exercise were asked to report the number of hours they spent per week on moderate (e.g., brisk walking) and vigorous (e.g., strenuous sports and jogging) exercise. Based on these data the researchers estimated the total hours of metabolic equivalent tasks (MET) per week, a value always greater than 0. The table below gives summary statistics of MET for women in this study based on the amount of coffee consumed.<sup>41</sup>

<i>Caffeinated coffee consumption</i>						
	$\leq 1$ cup/week	2-6 cups/week	1 cup/day	2-3 cups/day	$\geq 4$ cups/day	Total
Mean	18.7	19.6	19.3	18.9	17.5	
SD	21.1	25.5	22.5	22.0	22.0	
n	12,215	6,617	17,234	12,290	2,383	50,739

- (a) Write the hypotheses for evaluating if the average physical activity level varies among the different levels of coffee consumption.
- (b) Check conditions and describe any assumptions you must make to proceed with the test.
- (c) Below is part of the output associated with this test. Fill in the empty cells.

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
coffee	[ ]	[ ]	[ ]	[ ]	0.0003
Residuals	[ ]	25,564,819	[ ]		
Total	[ ]	25,575,327			

- (d) What is the conclusion of the test?

---

<sup>41</sup>M. Lucas et al. "Coffee, caffeine, and risk of depression among women". In: *Archives of internal medicine* 171.17 (2011), p. 1571.

**5.44 Student performance across discussion sections.** A professor who teaches a large introductory statistics class (197 students) with eight discussion sections would like to test if student performance differs by discussion section, where each discussion section has a different teaching assistant. The summary table below shows the average final exam score for each discussion section as well as the standard deviation of scores and the number of students in each section.

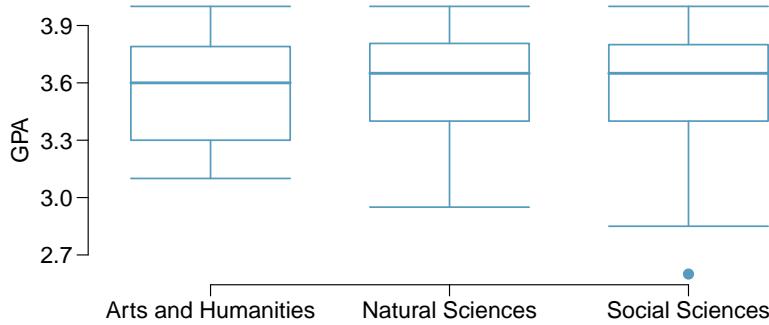
	Sec 1	Sec 2	Sec 3	Sec 4	Sec 5	Sec 6	Sec 7	Sec 8
$n_i$	33	19	10	29	33	10	32	31
$\bar{x}_i$	92.94	91.11	91.80	92.45	89.30	88.30	90.12	93.35
$s_i$	4.21	5.58	3.43	5.92	9.32	7.27	6.93	4.57

The ANOVA output below can be used to test for differences between the average scores from the different discussion sections.

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
section	7	525.01	75.00	1.87	0.0767
Residuals	189	7584.11	40.13		

Conduct a hypothesis test to determine if these data provide convincing evidence that the average score varies across some (or all) groups. Check conditions and describe any assumptions you must make to proceed with the test.

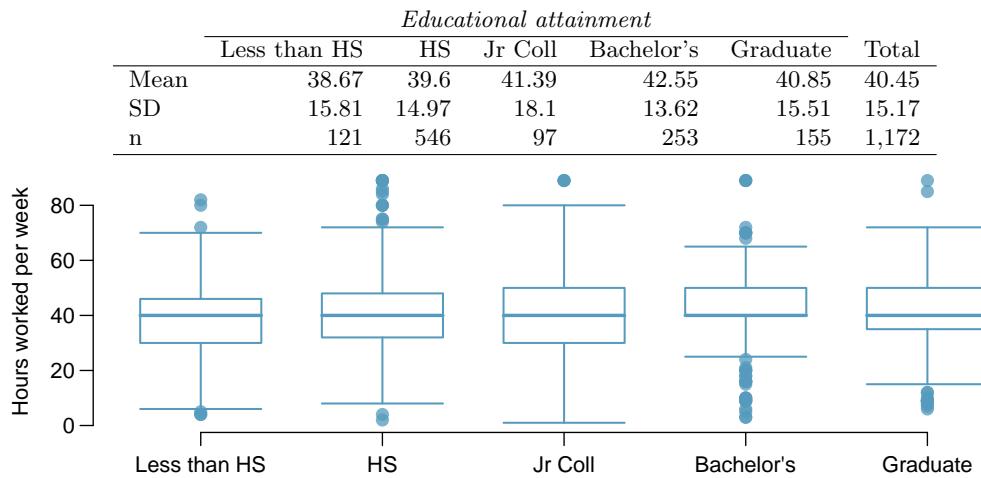
**5.45 GPA and major.** Undergraduate students taking an introductory statistics course at Duke University conducted a survey about GPA and major. The side-by-side box plots show the distribution of GPA among three groups of majors. Also provided is the ANOVA output.



	Df	Sum Sq	Mean Sq	F value	Pr(>F)
major	2	0.03	0.015	0.185	0.8313
Residuals	195	15.77	0.081		

- (a) Write the hypotheses for testing for a difference between average GPA across majors.
- (b) What is the conclusion of the hypothesis test?
- (c) How many students answered these questions on the survey, i.e. what is the sample size?

**5.46 Work hours and education.** The General Social Survey collects data on demographics, education, and work, among many other characteristics of US residents.<sup>42</sup> Using ANOVA, we can consider educational attainment levels for all 1,172 respondents at once. Below are the distributions of hours worked by educational attainment and relevant summary statistics that will be helpful in carrying out this analysis.



- (a) Write hypotheses for evaluating whether the average number of hours worked varies across the five groups.
- (b) Check conditions and describe any assumptions you must make to proceed with the test.
- (c) Below is part of the output associated with this test. Fill in the empty cells.

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
degree	[ ]	[ ]	501.54	[ ]	0.0682
Residuals	[ ]	267,382	[ ]		
Total	[ ]	[ ]			

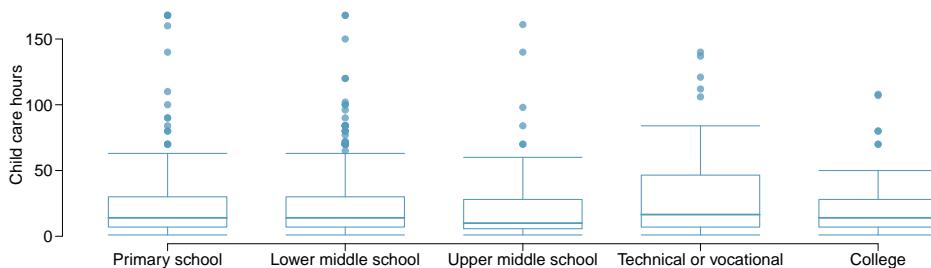
- (d) What is the conclusion of the test?

**5.47 True / False: ANOVA, Part I.** Determine if the following statements are true or false in ANOVA, and explain your reasoning for statements you identify as false.

- (a) As the number of groups increases, the modified significance level for pairwise tests increases as well.
- (b) As the total sample size increases, the degrees of freedom for the residuals increases as well.
- (c) The constant variance condition can be somewhat relaxed when the sample sizes are relatively consistent across groups.
- (d) The independence assumption can be relaxed when the total sample size is large.

<sup>42</sup>National Opinion Research Center, General Social Survey, 2010.

**5.48 Child care hours.** The China Health and Nutrition Survey aims to examine the effects of the health, nutrition, and family planning policies and programs implemented by national and local governments.<sup>43</sup> It, for example, collects information on number of hours Chinese parents spend taking care of their children under age 6. The side-by-side box plots below show the distribution of this variable by educational attainment of the parent. Also provided below is the ANOVA output for comparing average hours across educational attainment categories.



	Df	Sum Sq	Mean Sq	F value	Pr(>F)
education	4	4142.09	1035.52	1.26	0.2846
Residuals	794	653047.83	822.48		

- (a) Write the hypotheses for testing for a difference between the average number of hours spent on child care across educational attainment levels.
- (b) What is the conclusion of the hypothesis test?

**5.49 Prison isolation experiment, Part II.** Exercise 5.37 introduced an experiment that was conducted with the goal of identifying a treatment that reduces subjects' psychopathic deviant T scores, where this score measures a person's need for control or his rebellion against control. In Exercise 5.37 you evaluated the success of each treatment individually. An alternative analysis involves comparing the success of treatments. The relevant ANOVA output is given below.

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
treatment	2	639.48	319.74	3.33	0.0461
Residuals	39	3740.43	95.91		
$s_{pooled} = 9.793$ on $df = 39$					

- (a) What are the hypotheses?
- (b) What is the conclusion of the test? Use a 5% significance level.
- (c) If in part (b) you determined that the test is significant, conduct pairwise tests to determine which groups are different from each other. If you did not reject the null hypothesis in part (b), recheck your answer.

**5.50 True / False: ANOVA, Part II.** Determine if the following statements are true or false, and explain your reasoning for statements you identify as false.

If the null hypothesis that the means of four groups are all the same is rejected using ANOVA at a 5% significance level, then ...

- (a) we can then conclude that all the means are different from one another.
- (b) the standardized variability between groups is higher than the standardized variability within groups.
- (c) the pairwise analysis will identify at least one pair of means that are significantly different.
- (d) the appropriate  $\alpha$  to be used in pairwise comparisons is  $0.05 / 4 = 0.0125$  since there are four groups.

<sup>43</sup>UNC Carolina Population Center, China Health and Nutrition Survey, 2006.

## Chapter 7

# Introduction to linear regression

Linear regression is a very powerful statistical technique. Many people have some familiarity with regression just from reading the news, where graphs with straight lines are overlaid on scatterplots. Linear models can be used for prediction or to evaluate whether there is a linear relationship between two numerical variables.

Figure 7.1 shows two variables whose relationship can be modeled perfectly with a straight line. The equation for the line is

$$y = 5 + 57.49x$$

Imagine what a perfect linear relationship would mean: you would know the exact value of  $y$  just by knowing the value of  $x$ . This is unrealistic in almost any natural process. For example, if we took family income  $x$ , this value would provide some useful information about how much financial support  $y$  a college may offer a prospective student. However, there would still be variability in financial support, even when comparing students whose families have similar financial backgrounds.

Linear regression assumes that the relationship between two variables,  $x$  and  $y$ , can be modeled by a straight line:

$$y = \beta_0 + \beta_1 x \tag{7.1}$$

$\beta_0, \beta_1$   
Linear model parameters

where  $\beta_0$  and  $\beta_1$  represent two model parameters ( $\beta$  is the Greek letter *beta*). These parameters are estimated using data, and we write their point estimates as  $b_0$  and  $b_1$ . When we use  $x$  to predict  $y$ , we usually call  $x$  the explanatory or **predictor** variable, and we call  $y$  the response.

It is rare for all of the data to fall on a straight line, as seen in the three scatterplots in Figure 7.2. In each case, the data fall around a straight line, even if none of the observations fall exactly on the line. The first plot shows a relatively strong downward linear trend, where the remaining variability in the data around the line is minor relative to the strength of the relationship between  $x$  and  $y$ . The second plot shows an upward trend that, while evident, is not as strong as the first. The last plot shows a very weak downward trend in the data, so slight we can hardly notice it. In each of these examples, we will have some uncertainty regarding our estimates of the model parameters,  $\beta_0$  and  $\beta_1$ . For instance, we might wonder, should we move the line up or down a little, or should we tilt it more or less?

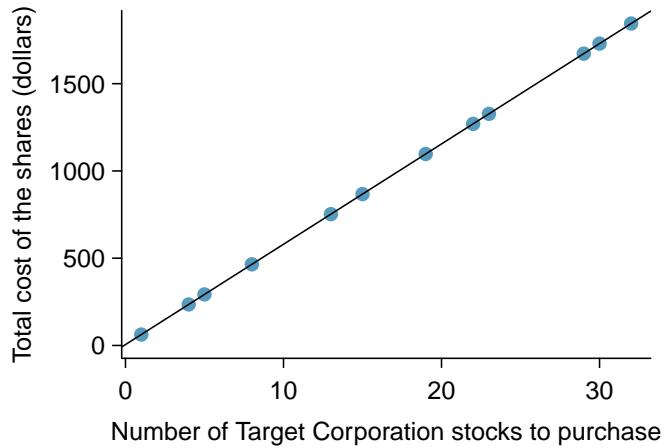


Figure 7.1: Requests from twelve separate buyers were simultaneously placed with a trading company to purchase Target Corporation stock (ticker TGT, April 26th, 2012), and the total cost of the shares were reported. Because the cost is computed using a linear formula, the linear fit is perfect.

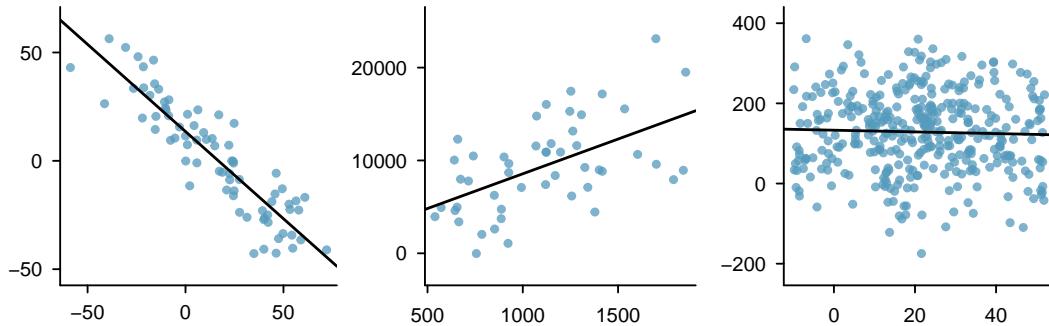


Figure 7.2: Three data sets where a linear model may be useful even though the data do not all fall exactly on the line.

As we move forward in this chapter, we will learn different criteria for line-fitting, and we will also learn about the uncertainty associated with estimates of model parameters.

We will also see examples in this chapter where fitting a straight line to the data, even if there is a clear relationship between the variables, is not helpful. One such case is shown in Figure 7.3 where there is a very strong relationship between the variables even though the trend is not linear. We will discuss nonlinear trends in this chapter and the next, but the details of fitting nonlinear models are saved for a later course.

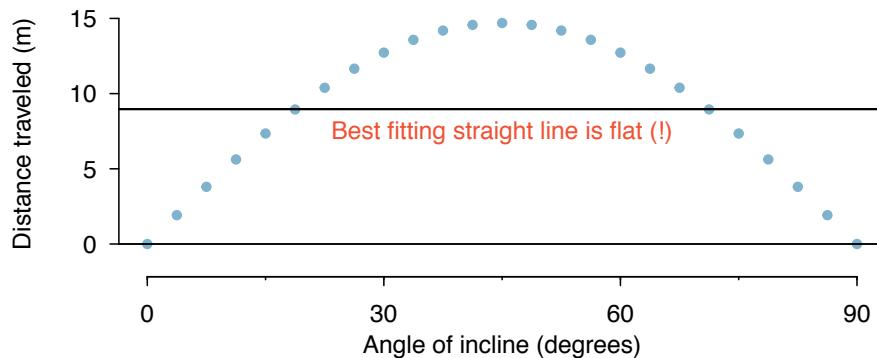


Figure 7.3: A linear model is not useful in this nonlinear case. These data are from an introductory physics experiment.

## 7.1 Line fitting, residuals, and correlation

It is helpful to think deeply about the line fitting process. In this section, we examine criteria for identifying a linear model and introduce a new statistic, *correlation*.

### 7.1.1 Beginning with straight lines

Scatterplots were introduced in Chapter 1 as a graphical technique to present two numerical variables simultaneously. Such plots permit the relationship between the variables to be examined with ease. Figure 7.4 shows a scatterplot for the head length and total length of 104 brushtail possums from Australia. Each point represents a single possum from the data.

The head and total length variables are associated. Possums with an above average total length also tend to have above average head lengths. While the relationship is not perfectly linear, it could be helpful to partially explain the connection between these variables with a straight line.

Straight lines should only be used when the data appear to have a linear relationship, such as the case shown in the left panel of Figure 7.6. The right panel of Figure 7.6 shows a case where a curved line would be more useful in understanding the relationship between the two variables.

#### Caution: Watch out for curved trends

We only consider models based on straight lines in this chapter. If data show a nonlinear trend, like that in the right panel of Figure 7.6, more advanced techniques should be used.

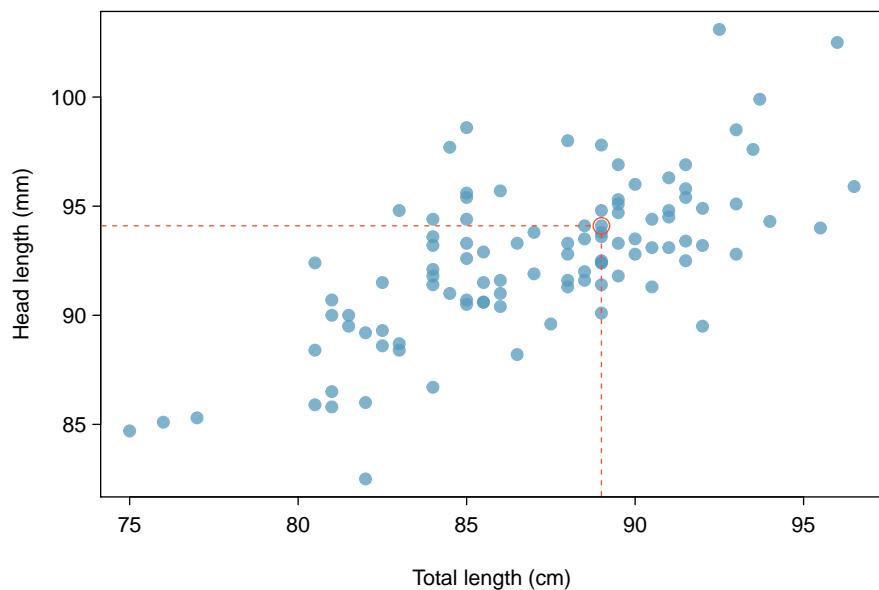


Figure 7.4: A scatterplot showing head length against total length for 104 brushtail possums. A point representing a possum with head length 94.1mm and total length 89cm is highlighted.



Figure 7.5: The common brushtail possum of Australia.

Photo by Greg Schechter (<https://flic.kr/p/9BAFbR>). CC BY 2.0 license.

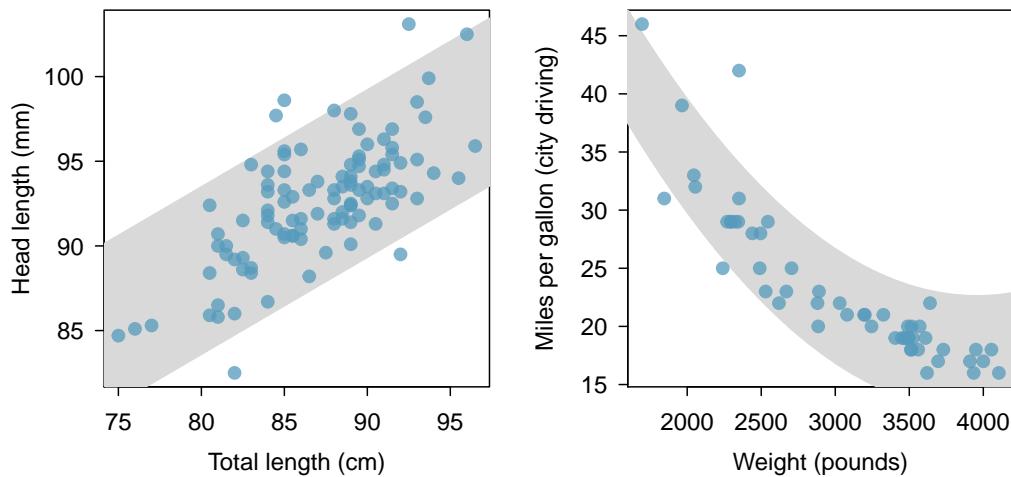


Figure 7.6: The figure on the left shows head length versus total length, and reveals that many of the points could be captured by a straight band. On the right, we see that a curved band is more appropriate in the scatterplot for `weight` and `mpgCity` from the `cars` data set.

### 7.1.2 Fitting a line by eye

We want to describe the relationship between the head length and total length variables in the possum data set using a line. In this example, we will use the total length as the predictor variable,  $x$ , to predict a possum's head length,  $y$ . We could fit the linear relationship by eye, as in Figure 7.7. The equation for this line is

$$\hat{y} = 41 + 0.59x \quad (7.2)$$

We can use this line to discuss properties of possums. For instance, the equation predicts a possum with a total length of 80 cm will have a head length of

$$\begin{aligned}\hat{y} &= 41 + 0.59 \times 80 \\ &= 88.2\end{aligned}$$

A “hat” on  $y$  is used to signify that this is an estimate. This estimate may be viewed as an average: the equation predicts that possums with a total length of 80 cm will have an average head length of 88.2 mm. Absent further information about an 80 cm possum, the prediction for head length that uses the average is a reasonable estimate.

### 7.1.3 Residuals

**Residuals** are the leftover variation in the data after accounting for the model fit:

$$\text{Data} = \text{Fit} + \text{Residual}$$

Each observation will have a residual. If an observation is above the regression line, then its residual, the vertical distance from the observation to the line, is positive. Observations below the line have negative residuals. One goal in picking the right linear model is for these residuals to be as small as possible.

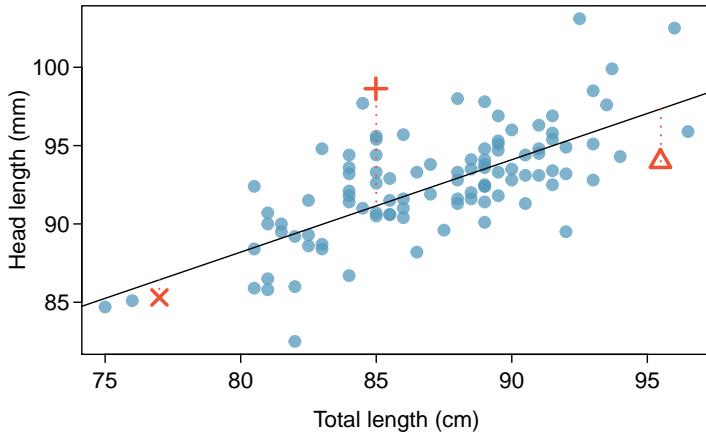


Figure 7.7: A reasonable linear model was fit to represent the relationship between head length and total length.

Three observations are noted specially in Figure 7.7. The observation marked by an “ $\times$ ” has a small, negative residual of about -1; the observation marked by “+” has a large residual of about +7; and the observation marked by “ $\triangle$ ” has a moderate residual of about -4. The size of a residual is usually discussed in terms of its absolute value. For example, the residual for “ $\triangle$ ” is larger than that of “ $\times$ ” because  $| - 4 |$  is larger than  $| - 1 |$ .

#### Residual: difference between observed and expected

The residual of the  $i^{th}$  observation  $(x_i, y_i)$  is the difference of the observed response  $(y_i)$  and the response we would predict based on the model fit  $(\hat{y}_i)$ :

$$e_i = y_i - \hat{y}_i$$

We typically identify  $\hat{y}_i$  by plugging  $x_i$  into the model.

- **Example 7.3** The linear fit shown in Figure 7.7 is given as  $\hat{y} = 41 + 0.59x$ . Based on this line, formally compute the residual of the observation  $(77.0, 85.3)$ . This observation is denoted by “ $\times$ ” on the plot. Check it against the earlier visual estimate, -1.

We first compute the predicted value of point “ $\times$ ” based on the model:

$$\hat{y}_{\times} = 41 + 0.59x_{\times} = 41 + 0.59 \times 77.0 = 86.4$$

Next we compute the difference of the actual head length and the predicted head length:

$$e_{\times} = y_{\times} - \hat{y}_{\times} = 85.3 - 86.4 = -1.1$$

This is very close to the visual estimate of -1.

- **Guided Practice 7.4** If a model underestimates an observation, will the residual be positive or negative? What about if it overestimates the observation?<sup>1</sup>

<sup>1</sup>If a model underestimates an observation, then the model estimate is below the actual. The residual, which is the actual observation value minus the model estimate, must then be positive. The opposite is true when the model overestimates the observation: the residual is negative.

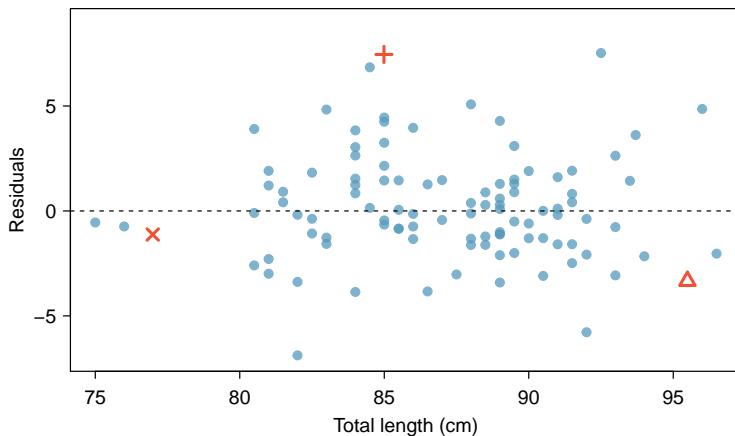


Figure 7.8: Residual plot for the model in Figure 7.7.

- Ⓐ **Guided Practice 7.5** Compute the residuals for the observations  $(85.0, 98.6)$  (“+” in the figure) and  $(95.5, 94.0)$  (“ $\Delta$ ”) using the linear relationship  $\hat{y} = 41 + 0.59x$ .<sup>2</sup>

Residuals are helpful in evaluating how well a linear model fits a data set. We often display them in a **residual plot** such as the one shown in Figure 7.8 for the regression line in Figure 7.7. The residuals are plotted at their original horizontal locations but with the vertical coordinate as the residual. For instance, the point  $(85.0, 98.6)_+$  had a residual of 7.45, so in the residual plot it is placed at  $(85.0, 7.45)$ . Creating a residual plot is sort of like tipping the scatterplot over so the regression line is horizontal.

- Ⓑ **Example 7.6** One purpose of residual plots is to identify characteristics or patterns still apparent in data after fitting a model. Figure 7.9 shows three scatterplots with linear models in the first row and residual plots in the second row. Can you identify any patterns remaining in the residuals?

In the first data set (first column), the residuals show no obvious patterns. The residuals appear to be scattered randomly around the dashed line that represents 0.

The second data set shows a pattern in the residuals. There is some curvature in the scatterplot, which is more obvious in the residual plot. We should not use a straight line to model these data. Instead, a more advanced technique should be used.

The last plot shows very little upwards trend, and the residuals also show no obvious patterns. It is reasonable to try to fit a linear model to the data. However, it is unclear whether there is statistically significant evidence that the slope parameter is different from zero. The point estimate of the slope parameter, labeled  $b_1$ , is not zero, but we might wonder if this could just be due to chance. We will address this sort of scenario in Section 7.4.

---

<sup>2</sup>(+) First compute the predicted value based on the model:

$$\hat{y}_+ = 41 + 0.59x_+ = 41 + 0.59 \times 85.0 = 91.15$$

Then the residual is given by

$$e_+ = y_+ - \hat{y}_+ = 98.6 - 91.15 = 7.45$$

This was close to the earlier estimate of 7.

(Δ)  $\hat{y}_\Delta = 41 + 0.59x_\Delta = 97.3$ .  $e_\Delta = y_\Delta - \hat{y}_\Delta = -3.3$ , close to the estimate of -4.

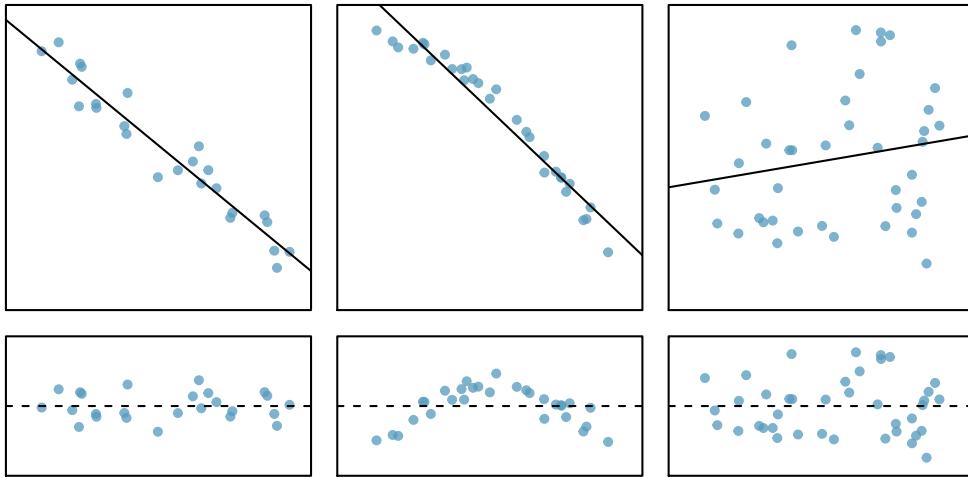


Figure 7.9: Sample data with their best fitting lines (top row) and their corresponding residual plots (bottom row).

### 7.1.4 Describing linear relationships with correlation

$R$   
correlation

#### Correlation: strength of a linear relationship

**Correlation**, which always takes values between -1 and 1, describes the strength of the linear relationship between two variables. We denote the correlation by  $R$ .

We can compute the correlation using a formula, just as we did with the sample mean and standard deviation. However, this formula is rather complex,<sup>3</sup> so we generally perform the calculations on a computer or calculator. Figure 7.10 shows eight plots and their corresponding correlations. Only when the relationship is perfectly linear is the correlation either -1 or 1. If the relationship is strong and positive, the correlation will be near +1. If it is strong and negative, it will be near -1. If there is no apparent linear relationship between the variables, then the correlation will be near zero.

The correlation is intended to quantify the strength of a linear trend. Nonlinear trends, even when strong, sometimes produce correlations that do not reflect the strength of the relationship; see three such examples in Figure 7.11.

**Guided Practice 7.7** It appears no straight line would fit any of the datasets represented in Figure 7.11. Try drawing nonlinear curves on each plot. Once you create a curve for each, describe what is important in your fit.<sup>4</sup>

<sup>3</sup>Formally, we can compute the correlation for observations  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  using the formula

$$R = \frac{1}{n-1} \sum_{i=1}^n \frac{x_i - \bar{x}}{s_x} \frac{y_i - \bar{y}}{s_y}$$

where  $\bar{x}$ ,  $\bar{y}$ ,  $s_x$ , and  $s_y$  are the sample means and standard deviations for each variable.

<sup>4</sup>We'll leave it to you to draw the lines. In general, the lines you draw should be close to most points and reflect overall trends in the data.

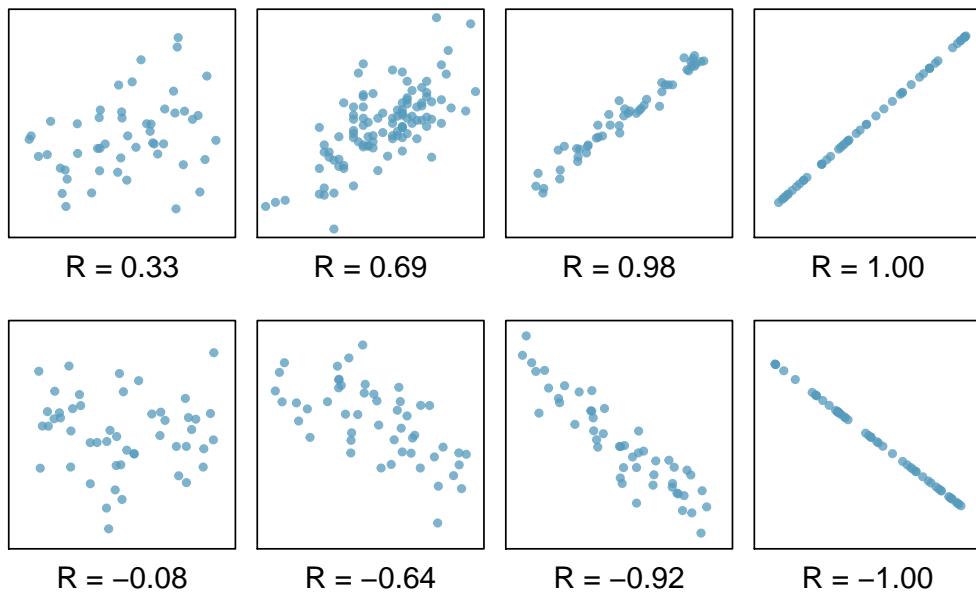


Figure 7.10: Sample scatterplots and their correlations. The first row shows variables with a positive relationship, represented by the trend up and to the right. The second row shows variables with a negative trend, where a large value in one variable is associated with a low value in the other.

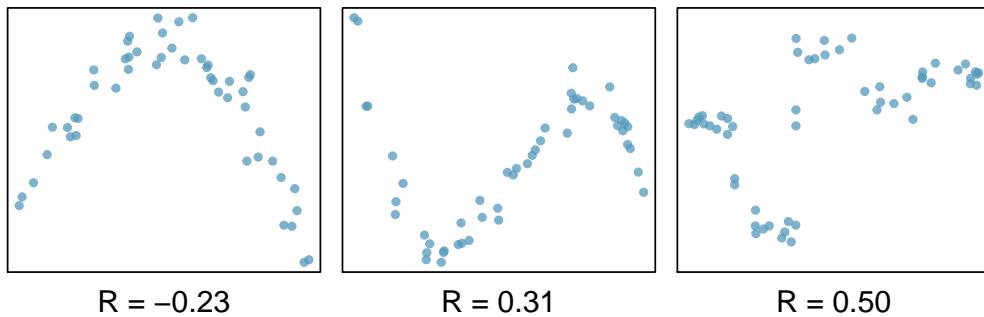


Figure 7.11: Sample scatterplots and their correlations. In each case, there is a strong relationship between the variables. However, the correlation is not very strong, and the relationship is not linear.

## 7.2 Fitting a line by least squares regression

Fitting linear models by eye is open to criticism since it is based on an individual preference. In this section, we use *least squares regression* as a more rigorous approach.

This section considers family income and gift aid data from a random sample of fifty students in the 2011 freshman class of Elmhurst College in Illinois.<sup>5</sup> Gift aid is financial aid that does not need to be paid back, as opposed to a loan. A scatterplot of the data is shown in Figure 7.12 along with two linear fits. The lines follow a negative trend in the data; students who have higher family incomes tended to have lower gift aid from the university.

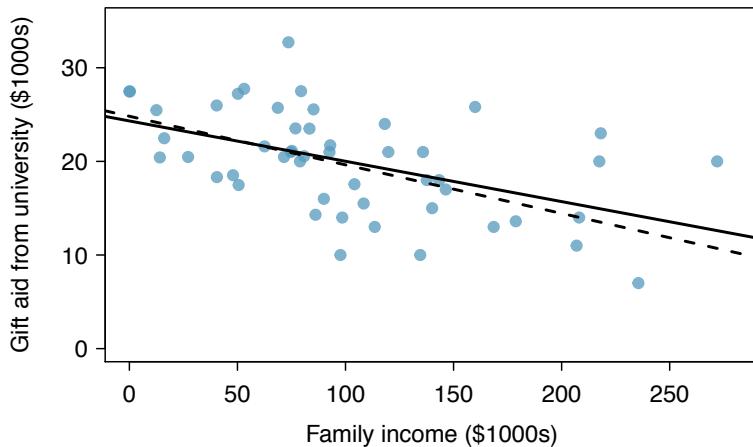


Figure 7.12: Gift aid and family income for a random sample of 50 freshman students from Elmhurst College. Two lines are fit to the data, the solid line being the *least squares line*.

 **Guided Practice 7.8** Is the correlation positive or negative in Figure 7.12?<sup>6</sup>

### 7.2.1 An objective measure for finding the best line

We begin by thinking about what we mean by “best”. Mathematically, we want a line that has small residuals. Perhaps our criterion could minimize the sum of the residual magnitudes:

$$|e_1| + |e_2| + \cdots + |e_n| \tag{7.9}$$

which we could accomplish with a computer program. The resulting dashed line shown in Figure 7.12 demonstrates this fit can be quite reasonable. However, a more common practice is to choose the line that minimizes the sum of the squared residuals:

$$e_1^2 + e_2^2 + \cdots + e_n^2 \tag{7.10}$$

<sup>5</sup>These data were sampled from a table of data for all freshman from the 2011 class at Elmhurst College that accompanied an article titled *What Students Really Pay to Go to College* published online by *The Chronicle of Higher Education*: [chronicle.com/article/What-Students-Really-Pay-to-Go/131435](http://chronicle.com/article/What-Students-Really-Pay-to-Go/131435)

<sup>6</sup>Larger family incomes are associated with lower amounts of aid, so the correlation will be negative. Using a computer, the correlation can be computed: -0.499.

The line that minimizes this **least squares criterion** is represented as the solid line in Figure 7.12. This is commonly called the **least squares line**. The following are three possible reasons to choose Criterion (7.10) over Criterion (7.9):

1. It is the most commonly used method.
2. Computing the line based on Criterion (7.10) is much easier by hand and in most statistical software.
3. In many applications, a residual twice as large as another residual is more than twice as bad. For example, being off by 4 is usually more than twice as bad as being off by 2. Squaring the residuals accounts for this discrepancy.

The first two reasons are largely for tradition and convenience; the last reason explains why Criterion (7.10) is typically most helpful.<sup>7</sup>

### 7.2.2 Conditions for the least squares line

When fitting a least squares line, we generally require

**Linearity.** The data should show a linear trend. If there is a nonlinear trend (e.g. left panel of Figure 7.13), an advanced regression method from another book or later course should be applied.

**Nearly normal residuals.** Generally the residuals must be nearly normal. When this condition is found to be unreasonable, it is usually because of outliers or concerns about influential points, which we will discuss in greater depth in Section 7.3. An example of non-normal residuals is shown in the second panel of Figure 7.13.

**Constant variability.** The variability of points around the least squares line remains roughly constant. An example of non-constant variability is shown in the third panel of Figure 7.13.

**Independent observations.** Be cautious about applying regression to **time series** data, which are sequential observations in time such as a stock price each day. Such data may have an underlying structure that should be considered in a model and analysis. An example of a data set where successive observations are not independent is shown in the fourth panel of Figure 7.13. There are also other instances where correlations within the data are important, which is further discussed in Chapter ??.

- ④ **Guided Practice 7.11** Should we have concerns about applying least squares regression to the Elmhurst data in Figure 7.12?<sup>8</sup>

---

<sup>7</sup>There are applications where Criterion (7.9) may be more useful, and there are plenty of other criteria we might consider. However, this book only applies the least squares criterion.

<sup>8</sup>The trend appears to be linear, the data fall around the line with no obvious outliers, the variance is roughly constant. These are also not time series observations. Least squares regression can be applied to these data.

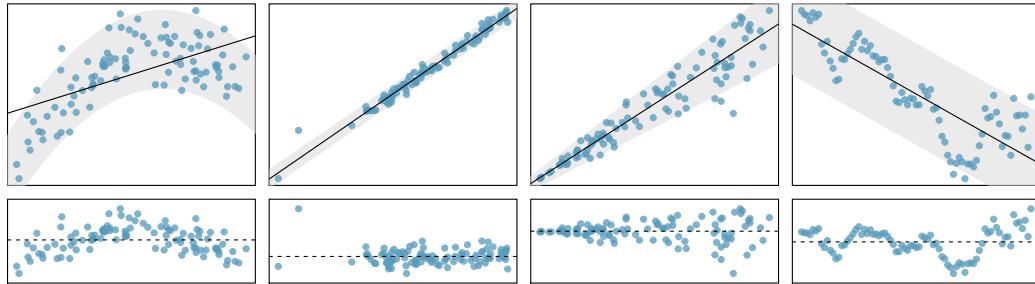


Figure 7.13: Four examples showing when the methods in this chapter are insufficient to apply to the data. In the left panel, a straight line does not fit the data. In the second panel, there are outliers; two points on the left are relatively distant from the rest of the data, and one of these points is very far away from the line. In the third panel, the variability of the data around the line increases with larger values of  $x$ . In the last panel, a time series data set is shown, where successive observations are highly correlated.

### 7.2.3 Finding the least squares line

For the Elmhurst data, we could write the equation of the least squares regression line as

$$\widehat{\text{aid}} = \beta_0 + \beta_1 \times \text{family\_income}$$

Here the equation is set up to predict gift aid based on a student's family income, which would be useful to students considering Elmhurst. These two values,  $\beta_0$  and  $\beta_1$ , are the *parameters* of the regression line.

As in Chapters 4–6, the parameters are estimated using observed data. In practice, this estimation is done using a computer in the same way that other estimates, like a sample mean, can be estimated using a computer or calculator. However, we can also find the parameter estimates by applying two properties of the least squares line:

- The slope of the least squares line can be estimated by

$$b_1 = \frac{s_y}{s_x} R \quad (7.12)$$

where  $R$  is the correlation between the two variables, and  $s_x$  and  $s_y$  are the sample standard deviations of the explanatory variable and response, respectively.

- If  $\bar{x}$  is the mean of the horizontal variable (from the data) and  $\bar{y}$  is the mean of the vertical variable, then the point  $(\bar{x}, \bar{y})$  is on the least squares line.

We use  $b_0$  and  $b_1$  to represent the point estimates of the parameters  $\beta_0$  and  $\beta_1$ .

-  **Guided Practice 7.13** Table 7.14 shows the sample means for the family income and gift aid as \$101,800 and \$19,940, respectively. Plot the point  $(101.8, 19.94)$  on Figure 7.12 on page 198 to verify it falls on the least squares line (the solid line).<sup>9</sup>

---

<sup>9</sup>If you need help finding this location, draw a straight line up from the  $x$ -value of 100 (or thereabout). Then draw a horizontal line at 20 (or thereabout). These lines should intersect on the least squares line.

	family income, in \$1000s ("x")	gift aid, in \$1000s ("y")
mean	$\bar{x} = 101.8$	$\bar{y} = 19.94$
sd	$s_x = 63.2$	$s_y = 5.46$
		$R = -0.499$

Table 7.14: Summary statistics for family income and gift aid.

- Ⓐ **Guided Practice 7.14** Using the summary statistics in Table 7.14, compute the slope for the regression line of gift aid against family income.<sup>10</sup>

You might recall the **point-slope** form of a line from math class (another common form is *slope-intercept*). Given the slope of a line and a point on the line,  $(x_0, y_0)$ , the equation for the line can be written as

$$y - y_0 = \text{slope} \times (x - x_0) \quad (7.15)$$

A common exercise to become more familiar with foundations of least squares regression is to use basic summary statistics and point-slope form to produce the least squares line.

**TIP: Identifying the least squares line from summary statistics**

To identify the least squares line from summary statistics:

- Estimate the slope parameter,  $b_1$ , using Equation (7.12).
- Noting that the point  $(\bar{x}, \bar{y})$  is on the least squares line, use  $x_0 = \bar{x}$  and  $y_0 = \bar{y}$  along with the slope  $b_1$  in the point-slope equation:

$$y - \bar{y} = b_1(x - \bar{x})$$

- Simplify the equation.

- Ⓑ **Example 7.16** Using the point  $(101.8, 19.94)$  from the sample means and the slope estimate  $b_1 = -0.0431$  from Guided Practice 7.14, find the least-squares line for predicting aid based on family income.

Apply the point-slope equation using  $(101.8, 19.94)$  and the slope  $b_1 = -0.0431$ :

$$\begin{aligned} y - y_0 &= b_1(x - x_0) \\ y - 19.94 &= -0.0431(x - 101.8) \end{aligned}$$

Expanding the right side and then adding 19.94 to each side, the equation simplifies:

$$\widehat{\text{aid}} = 24.3 - 0.0431 \times \text{family\_income}$$

Here we have replaced  $y$  with  $\widehat{\text{aid}}$  and  $x$  with  $\text{family\_income}$  to put the equation in context.

---

<sup>10</sup>Apply Equation (7.12) with the summary statistics from Table 7.14 to compute the slope:

$$b_1 = \frac{s_y}{s_x} R = \frac{5.46}{63.2} (-0.499) = -0.0431$$

We mentioned earlier that a computer is usually used to compute the least squares line. A summary table based on computer output is shown in Table 7.15 for the Elmhurst data. The first column of numbers provides estimates for  $b_0$  and  $b_1$ , respectively. Compare these to the result from Example 7.16.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	24.3193	1.2915	18.83	0.0000
family_income	-0.0431	0.0108	-3.98	0.0002

Table 7.15: Summary of least squares fit for the Elmhurst data. Compare the parameter estimates in the first column to the results of Example 7.16.

- **Example 7.17** Examine the second, third, and fourth columns in Table 7.15. Can you guess what they represent?

---

We'll describe the meaning of the columns using the second row, which corresponds to  $\beta_1$ . The first column provides the point estimate for  $\beta_1$ , as we calculated in an earlier example: -0.0431. The second column is a standard error for this point estimate: 0.0108. The third column is a  $t$ -test statistic for the null hypothesis that  $\beta_1 = 0$ :  $T = -3.98$ . The last column is the p-value for the  $t$ -test statistic for the null hypothesis  $\beta_1 = 0$  and a two-sided alternative hypothesis: 0.0002. We will get into more of these details in Section 7.4.

- **Example 7.18** Suppose a high school senior is considering Elmhurst College. Can she simply use the linear equation that we have estimated to calculate her financial aid from the university?

---

She may use it as an estimate, though some qualifiers on this approach are important. First, the data all come from one freshman class, and the way aid is determined by the university may change from year to year. Second, the equation will provide an imperfect estimate. While the linear equation is good at capturing the trend in the data, no individual student's aid will be perfectly predicted.

### 7.2.4 Interpreting regression line parameter estimates

Interpreting parameters in a regression model is often one of the most important steps in the analysis.

- **Example 7.19** The slope and intercept estimates for the Elmhurst data are -0.0431 and 24.3. What do these numbers really mean?

Interpreting the slope parameter is helpful in almost any application. For each additional \$1,000 of family income, we would expect a student to receive a net difference of  $\$1,000 \times (-0.0431) = -\$43.10$  in aid on average, i.e. *\$43.10 less*. Note that a higher family income corresponds to less aid because the coefficient of family income is negative in the model. We must be cautious in this interpretation: while there is a real association, we cannot interpret a causal connection between the variables because these data are observational. That is, increasing a student's family income may not cause the student's aid to drop. (It would be reasonable to contact the college and ask if the relationship is causal, i.e. if Elmhurst College's aid decisions are partially based on students' family income.)

The estimated intercept  $b_0 = 24.3$  (in \$1000s) describes the average aid if a student's family had no income. The meaning of the intercept is relevant to this application since the family income for some students at Elmhurst is \$0. In other applications, the intercept may have little or no practical value if there are no observations where  $x$  is near zero.

#### Interpreting parameters estimated by least squares

The slope describes the estimated difference in the  $y$  variable if the explanatory variable  $x$  for a case happened to be one unit larger. The intercept describes the average outcome of  $y$  if  $x = 0$  and the linear model is valid all the way to  $x = 0$ , which in many applications is not the case.

### 7.2.5 Extrapolation is treacherous

*When those blizzards hit the East Coast this winter, it proved to my satisfaction that global warming was a fraud. That snow was freezing cold. But in an alarming trend, temperatures this spring have risen. Consider this: On February 6<sup>th</sup> it was 10 degrees. Today it hit almost 80. At this rate, by August it will be 220 degrees. So clearly folks the climate debate rages on.*

Stephen Colbert  
April 6th, 2010<sup>11</sup>

Linear models can be used to approximate the relationship between two variables. However, these models have real limitations. Linear regression is simply a modeling framework. The truth is almost always much more complex than our simple line. For example, we do not know how the data outside of our limited window will behave.

<sup>11</sup>[www.cc.com/video-clips/l4nkoq](http://www.cc.com/video-clips/l4nkoq)

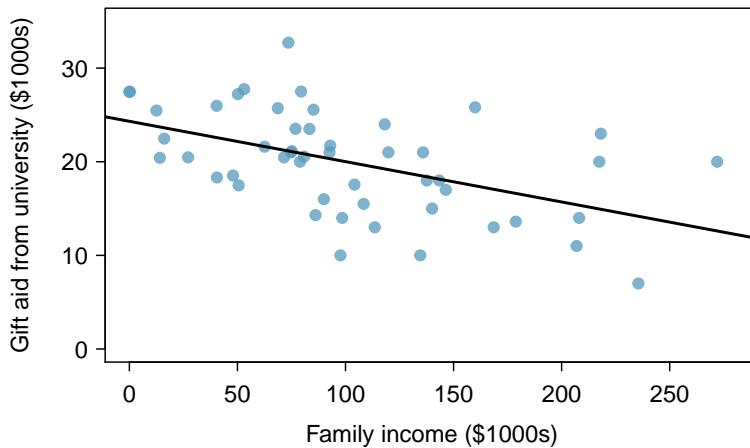


Figure 7.16: Gift aid and family income for a random sample of 50 freshman students from Elmhurst College, shown with the least squares regression line.

- **Example 7.20** Use the model  $\widehat{aid} = 24.3 - 0.0431 \times \text{family\_income}$  to estimate the aid of another freshman student whose family had income of \$1 million.

---

Recall that the units of family income are in \$1000s, so we want to calculate the aid for  $\text{family\_income} = 1000$ :

$$24.3 - 0.0431 \times \text{family\_income} = 24.3 - 0.0431 \times 1000 = -18.8$$

The model predicts this student will have -\$18,800 in aid (!). Elmhurst College cannot (or at least does not) require any students to pay extra on top of tuition to attend.

Applying a model estimate to values outside of the realm of the original data is called **extrapolation**. Generally, a linear model is only an approximation of the real relationship between two variables. If we extrapolate, we are making an unreliable bet that the approximate linear relationship will be valid in places where it has not been analyzed.

### 7.2.6 Using $R^2$ to describe the strength of a fit

We evaluated the strength of the linear relationship between two variables earlier using the correlation,  $R$ . However, it is more common to explain the strength of a linear fit using  $R^2$ , called **R-squared**. If provided with a linear model, we might like to describe how closely the data cluster around the linear fit.

The  $R^2$  of a linear model describes the amount of variation in the response that is explained by the least squares line. For example, consider the Elmhurst data, shown in Figure 7.16. The variance of the response variable, aid received, is  $s_{\text{aid}}^2 = 29.8$ . However, if we apply our least squares line, then this model reduces our uncertainty in predicting aid using a student's family income. The variability in the residuals describes how much variation remains after using the model:  $s_{\text{RES}}^2 = 22.4$ . In short, there was a reduction of

$$\frac{s_{\text{aid}}^2 - s_{\text{RES}}^2}{s_{\text{aid}}^2} = \frac{29.8 - 22.4}{29.8} = \frac{7.5}{29.8} = 0.25$$

or about 25% in the data's variation by using information about family income for predicting aid using a linear model. This corresponds exactly to the R-squared value:

$$R = -0.499$$

$$R^2 = 0.25$$

- Ⓐ **Guided Practice 7.21** If a linear model has a very strong negative relationship with a correlation of -0.97, how much of the variation in the response is explained by the explanatory variable?<sup>12</sup>

#### Calculator videos

Videos covering how to find regression coefficients using TI and Casio graphing calculators are available at [openintro.org/videos](http://openintro.org/videos).

### 7.2.7 Categorical predictors with two levels

Categorical variables are also useful in predicting outcomes. Here we consider a categorical predictor with two levels (recall that a *level* is the same as a *category*). We'll consider Ebay auctions for a video game, *Mario Kart* for the Nintendo Wii, where both the total price of the auction and the condition of the game were recorded.<sup>13</sup> Here we want to predict total price based on game condition, which takes values `used` and `new`. A plot of the auction data is shown in Figure 7.17.

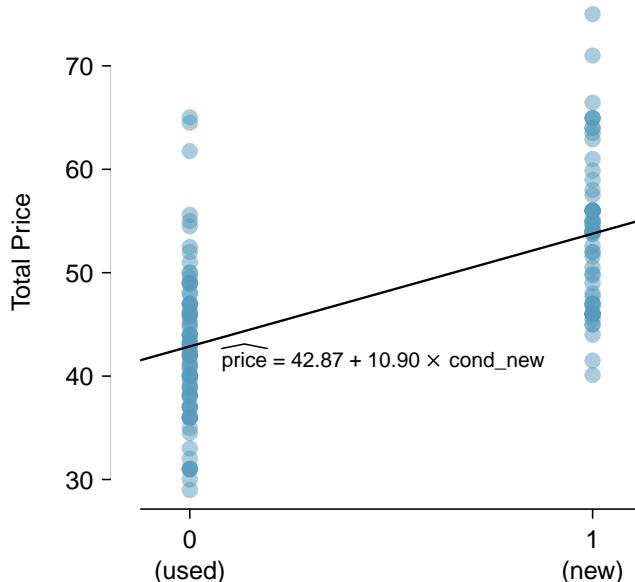


Figure 7.17: Total auction prices for the video game *Mario Kart*, divided into used ( $x = 0$ ) and new ( $x = 1$ ) condition games. The least squares regression line is also shown.

<sup>12</sup>About  $R^2 = (-0.97)^2 = 0.94$  or 94% of the variation is explained by the linear model.

<sup>13</sup>These data were collected in Fall 2009 and may be found at [openintro.org](http://openintro.org).

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	42.87	0.81	52.67	0.0000
cond_new	10.90	1.26	8.66	0.0000

Table 7.18: Least squares regression summary for the final auction price against the condition of the game.

To incorporate the game condition variable into a regression equation, we must convert the categories into a numerical form. We will do so using an **indicator variable** called `cond_new`, which takes value 1 when the game is new and 0 when the game is used. Using this indicator variable, the linear model may be written as

$$\widehat{\text{price}} = \beta_0 + \beta_1 \times \text{cond\_new}$$

The fitted model is summarized in Table 7.18, and the model with its parameter estimates is given as

$$\widehat{\text{price}} = 42.87 + 10.90 \times \text{cond\_new}$$

For categorical predictors with just two levels, the linearity assumption will always be satisfied. However, we must evaluate whether the residuals in each group are approximately normal and have approximately equal variance. As can be seen in Figure 7.17, both of these conditions are reasonably satisfied by the auction data.

- **Example 7.22** Interpret the two parameters estimated in the model for the price of *Mario Kart* in eBay auctions.

---

The intercept is the estimated price when `cond_new` takes value 0, i.e. when the game is in used condition. That is, the average selling price of a used version of the game is \$42.87.

The slope indicates that, on average, new games sell for about \$10.90 more than used games.

**TIP: Interpreting model estimates for categorical predictors.**

The estimated intercept is the value of the response variable for the first category (i.e. the category corresponding to an indicator value of 0). The estimated slope is the average change in the response variable between the two categories.

We'll elaborate further on this Ebay auction data in Chapter ??, where we examine the influence of many predictor variables simultaneously using multiple regression. In multiple regression, we will consider the association of auction price with regard to each variable while controlling for the influence of other variables. This is especially important since some of the predictors are associated. For example, auctions with games in new condition also often came with more accessories.

## 7.3 Types of outliers in linear regression



In this section, we identify criteria for determining which outliers are important and influential. Outliers in regression are observations that fall far from the “cloud” of points. These points are especially important because they can have a strong influence on the least squares line.

**Example 7.23** There are six plots shown in Figure 7.19 along with the least squares line and residual plots. For each scatterplot and residual plot pair, identify the outliers and note how they influence the least squares line. Recall that an outlier is any point that doesn’t appear to belong with the vast majority of the other points.

- (1) There is one outlier far from the other points, though it only appears to slightly influence the line.
- (2) There is one outlier on the right, though it is quite close to the least squares line, which suggests it wasn’t very influential.
- (3) There is one point far away from the cloud, and this outlier appears to pull the least squares line up on the right; examine how the line around the primary cloud doesn’t appear to fit very well.
- (4) There is a primary cloud and then a small secondary cloud of four outliers. The secondary cloud appears to be influencing the line somewhat strongly, making the least square line fit poorly almost everywhere. There might be an interesting explanation for the dual clouds, which is something that could be investigated.
- (5) There is no obvious trend in the main cloud of points and the outlier on the right appears to largely control the slope of the least squares line.
- (6) There is one outlier far from the cloud, however, it falls quite close to the least squares line and does not appear to be very influential.

Examine the residual plots in Figure 7.19. You will probably find that there is some trend in the main clouds of (3) and (4). In these cases, the outliers influenced the slope of the least squares lines. In (5), data with no clear trend were assigned a line with a large trend simply due to one outlier (!).

### Leverage

Points that fall horizontally away from the center of the cloud tend to pull harder on the line, so we call them points with **high leverage**.

Points that fall horizontally far from the line are points of high leverage; these points can strongly influence the slope of the least squares line. If one of these high leverage points does appear to actually invoke its influence on the slope of the line – as in cases (3), (4), and (5) of Example 7.23 – then we call it an **influential point**. Usually we can say a point is influential if, had we fitted the line without it, the influential point would have been unusually far from the least squares line.

It is tempting to remove outliers. Don’t do this without a very good reason. Models that ignore exceptional (and interesting) cases often perform poorly. For instance, if a financial firm ignored the largest market swings – the “outliers” – they would soon go bankrupt by making poorly thought-out investments.

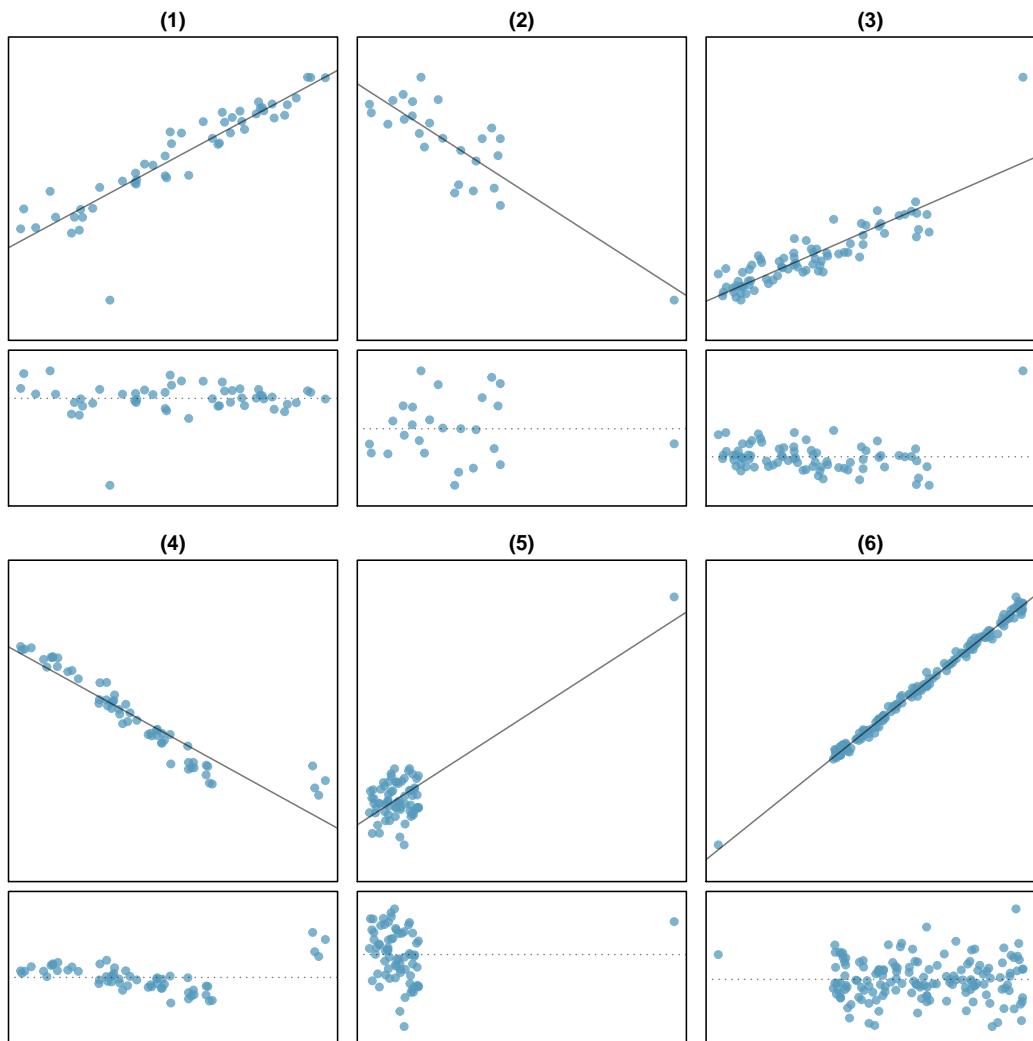


Figure 7.19: Six plots, each with a least squares line and residual plot. All data sets have at least one outlier.

**Caution: Don't ignore outliers when fitting a final model**

If there are outliers in the data, they should not be removed or ignored without a good reason. Whatever final model is fit to the data would not be very helpful if it ignores the most exceptional cases.

**Caution: Outliers for a categorical predictor with two levels**

Be cautious about using a categorical predictor when one of the levels has very few observations. When this happens, those few observations become influential points.

## 7.4 Inference for linear regression

In this section we discuss uncertainty in the estimates of the slope and y-intercept for a regression line. Just as we identified standard errors for point estimates in previous chapters, we first discuss standard errors for these new estimates. However, in the case of regression, we will identify standard errors using statistical software.

### 7.4.1 Midterm elections and unemployment

Elections for members of the United States House of Representatives occur every two years, coinciding every four years with the U.S. Presidential election. The set of House elections occurring during the middle of a Presidential term are called midterm elections. In America's two-party system, one political theory suggests the higher the unemployment rate, the worse the President's party will do in the midterm elections.

To assess the validity of this claim, we can compile historical data and look for a connection. We consider every midterm election from 1898 to 2010, with the exception of those elections during the Great Depression. Figure 7.20 shows these data and the least-squares regression line:

$$\begin{aligned} \text{\% change in House seats for President's party} \\ = -6.71 - 1.00 \times (\text{unemployment rate}) \end{aligned}$$

We consider the percent change in the number of seats of the President's party (e.g. percent change in the number of seats for Democrats in 2010) against the unemployment rate.

Examining the data, there are no clear deviations from linearity, the constant variance condition, or in the normality of residuals (though we don't examine a normal probability plot here). While the data are collected sequentially, a separate analysis was used to check for any apparent correlation between successive observations; no such correlation was found.

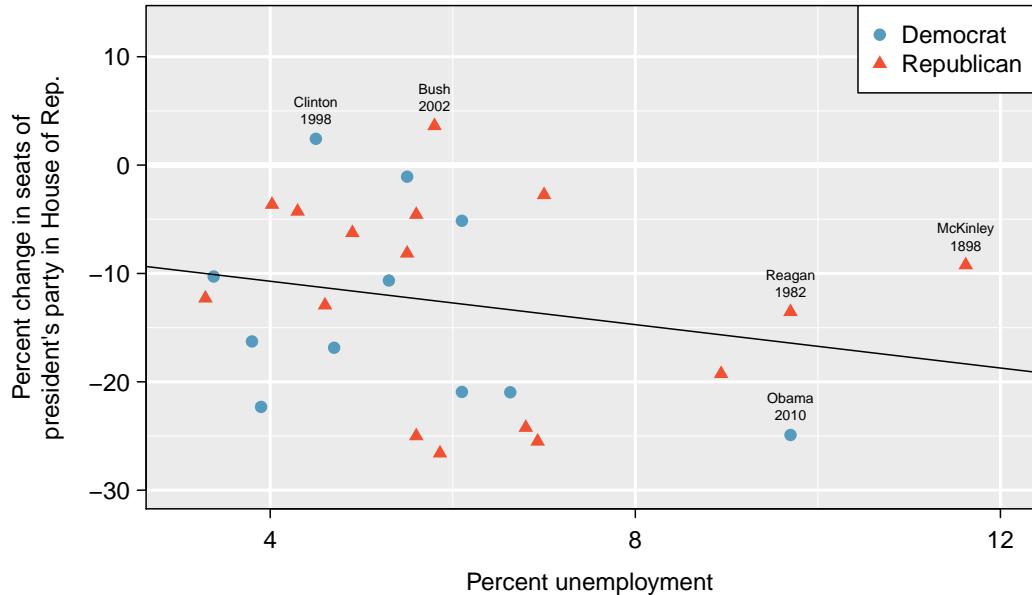


Figure 7.20: The percent change in House seats for the President’s party in each election from 1898 to 2010 plotted against the unemployment rate. The two points for the Great Depression have been removed, and a least squares regression line has been fit to the data.

Ⓐ **Guided Practice 7.24** The data for the Great Depression (1934 and 1938) were removed because the unemployment rate was 21% and 18%, respectively. Do you agree that they should be removed for this investigation? Why or why not?<sup>14</sup>

There is a negative slope in the line shown in Figure 7.20. However, this slope (and the y-intercept) are only estimates of the parameter values. We might wonder, is this convincing evidence that the “true” linear model has a negative slope? That is, do the data provide strong evidence that the political theory is accurate? We can frame this investigation into a one-sided statistical hypothesis test:

$$H_0: \beta_1 = 0. \text{ The true linear model has slope zero.}$$

$$H_A: \beta_1 < 0. \text{ The true linear model has a slope less than zero. The higher the unemployment, the greater the loss for the President’s party in the House of Representatives.}$$

We would reject  $H_0$  in favor of  $H_A$  if the data provide strong evidence that the true slope parameter is less than zero. To assess the hypotheses, we identify a standard error for the estimate, compute an appropriate test statistic, and identify the p-value.

<sup>14</sup>We will provide two considerations. Each of these points would have very high leverage on any least-squares regression line, and years with such high unemployment may not help us understand what would happen in other years where the unemployment is only modestly high. On the other hand, these are exceptional cases, and we would be discarding important information if we exclude them from a final analysis.

### 7.4.2 Understanding regression output from software

Just like other point estimates we have seen before, we can compute a standard error and test statistic for  $b_1$ . We will generally label the test statistic using a  $T$ , since it follows the  $t$ -distribution.

We will rely on statistical software to compute the standard error and leave the explanation of how this standard error is determined to a second or third statistics course. Table 7.21 shows software output for the least squares regression line in Figure 7.20. The row labeled *unemp* represents the information for the slope, which is the coefficient of the unemployment variable.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-6.7142	5.4567	-1.23	0.2300
unemp	-1.0010	0.8717	-1.15	0.2617
<i>df</i> = 25				

Table 7.21: Output from statistical software for the regression line modeling the midterm election losses for the President's party as a response to unemployment.

- Example 7.25 What do the first and second columns of Table 7.21 represent?

The entries in the first column represent the least squares estimates,  $b_0$  and  $b_1$ , and the values in the second column correspond to the standard errors of each estimate.

We previously used a  $t$ -test statistic for hypothesis testing in the context of numerical data. Regression is very similar. In the hypotheses we consider, the null value for the slope is 0, so we can compute the test statistic using the T (or Z) score formula:

$$T = \frac{\text{estimate} - \text{null value}}{\text{SE}} = \frac{-1.0010 - 0}{0.8717} = -1.15$$

We can look for the one-sided p-value – shown in Figure 7.22 – using the probability table for the  $t$ -distribution in Appendix B.2 on page 255.

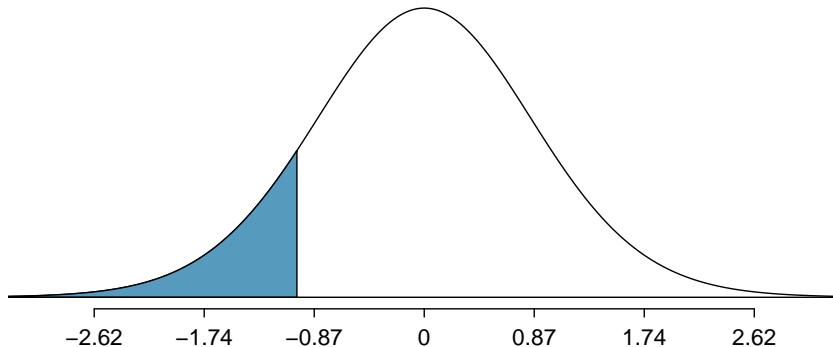


Figure 7.22: The distribution shown here is the sampling distribution for  $b_1$ , if the null hypothesis was true. The shaded tail represents the p-value for the hypothesis test evaluating whether there is convincing evidence that higher unemployment corresponds to a greater loss of House seats for the President's party during a midterm election.

- **Example 7.26** Table 7.21 offers the degrees of freedom for the test statistic  $T$ :  $df = 25$ . Identify the p-value for the hypothesis test.

Looking in the 25 degrees of freedom row in Appendix B.2, we see that the absolute value of the test statistic is smaller than any value listed, which means the tail area and therefore also the p-value is larger than 0.100 (one tail!). Because the p-value is so large, we fail to reject the null hypothesis. That is, the data do not provide convincing evidence that a higher unemployment rate has any correspondence with smaller or larger losses for the President's party in the House of Representatives in midterm elections.

We could have identified the  $t$ -test statistic from the software output in Table 7.21, shown in the second row (unemp) and third column (t value). The entry in the second row and last column in Table 7.21 represents the p-value for the two-sided hypothesis test where the null value is zero. The corresponding one-sided test would have a p-value half of the listed value.

#### Inference for regression

We usually rely on statistical software to identify point estimates and standard errors for parameters of a regression line. After verifying conditions hold for fitting a line, we can use the methods learned in Section 5.1 for the  $t$ -distribution to create confidence intervals for regression parameters or to evaluate hypothesis tests.

#### Caution: Don't carelessly use the p-value from regression output

The last column in regression output often lists p-values for one particular hypothesis: a two-sided test where the null value is zero. If your test is one-sided and the point estimate is in the direction of  $H_A$ , then you can halve the software's p-value to get the one-tail area. If neither of these scenarios match your hypothesis test, be cautious about using the software output to obtain the p-value.

- **Example 7.27** Examine Figure 7.16 on page 204, which relates the Elmhurst College aid and student family income. How sure are you that the slope is statistically significantly different from zero? That is, do you think a formal hypothesis test would reject the claim that the true slope of the line should be zero?

While the relationship between the variables is not perfect, there is an evident decreasing trend in the data. This suggests the hypothesis test will reject the null claim that the slope is zero.

- **Guided Practice 7.28** Table 7.23 shows statistical software output from fitting the least squares regression line shown in Figure 7.16. Use this output to formally evaluate the following hypotheses.  $H_0$ : The true coefficient for family income is zero.  $H_A$ : The true coefficient for family income is not zero.<sup>15</sup>

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	24.3193	1.2915	18.83	0.0000
family_income	-0.0431	0.0108	-3.98	0.0002
<i>df</i> = 48				

Table 7.23: Summary of least squares fit for the Elmhurst College data.

**TIP: Always check assumptions**

If conditions for fitting the regression line do not hold, then the methods presented here should not be applied. The standard error or distribution assumption of the point estimate – assumed to be normal when applying the *t*-test statistic – may not be valid.

 **Calculator videos**

Videos covering hypothesis testing for a regression coefficient using TI and Casio graphing calculators are available at [openintro.org/videos](http://openintro.org/videos).

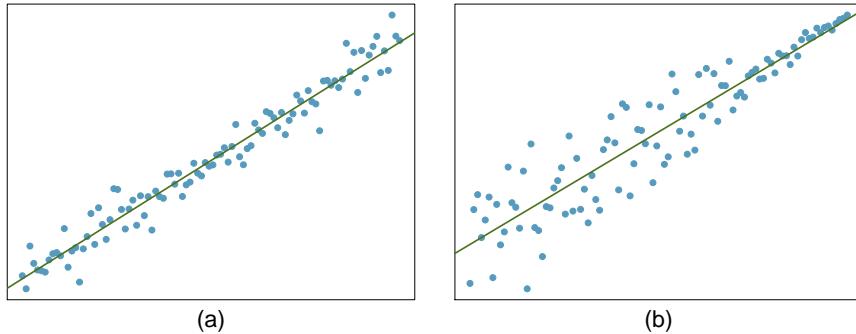
---

<sup>15</sup>We look in the second row corresponding to the family income variable. We see the point estimate of the slope of the line is -0.0431, the standard error of this estimate is 0.0108, and the *t*-test statistic is -3.98. The p-value corresponds exactly to the two-sided test we are interested in: 0.0002. The p-value is so small that we reject the null hypothesis and conclude that family income and financial aid at Elmhurst College for freshman entering in the year 2011 are negatively correlated and the true slope parameter is indeed less than 0, just as we believed in Example 7.27.

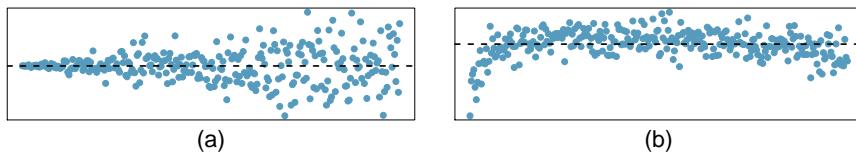
## 7.5 Exercises

### 7.5.1 Line fitting, residuals, and correlation

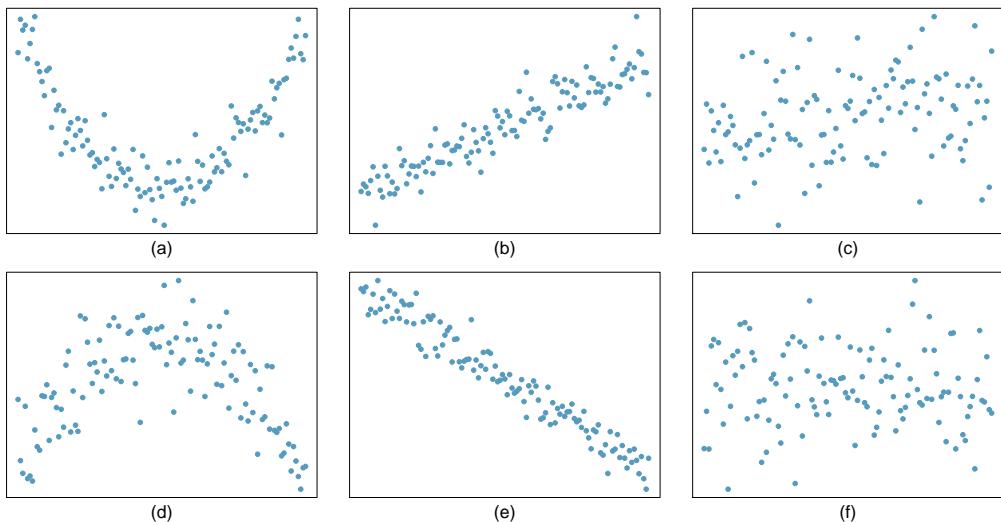
**7.1 Visualize the residuals.** The scatterplots shown below each have a superimposed regression line. If we were to construct a residual plot (residuals versus  $x$ ) for each, describe what those plots would look like.



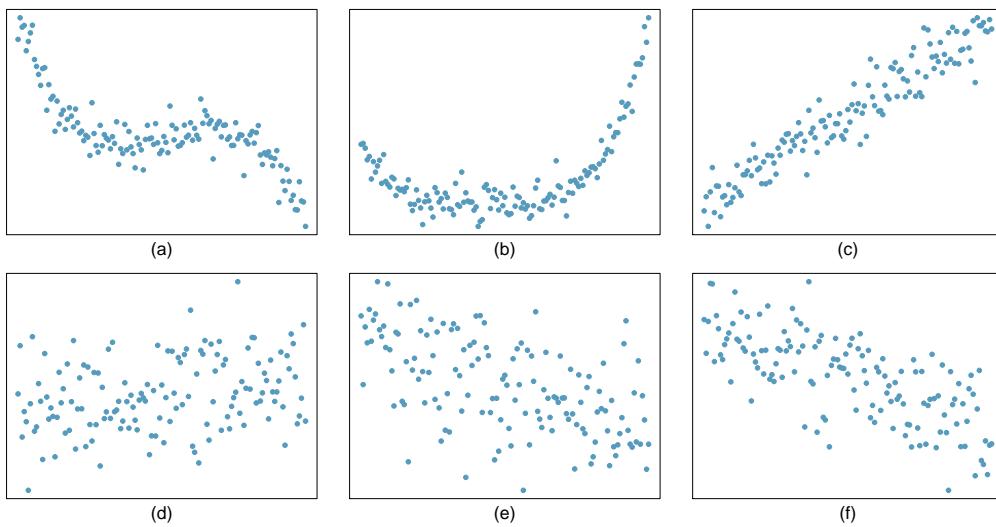
**7.2 Trends in the residuals.** Shown below are two plots of residuals remaining after fitting a linear model to two different sets of data. Describe important features and determine if a linear model would be appropriate for these data. Explain your reasoning.



**7.3 Identify relationships, Part I.** For each of the six plots, identify the strength of the relationship (e.g. weak, moderate, or strong) in the data and whether fitting a linear model would be reasonable.

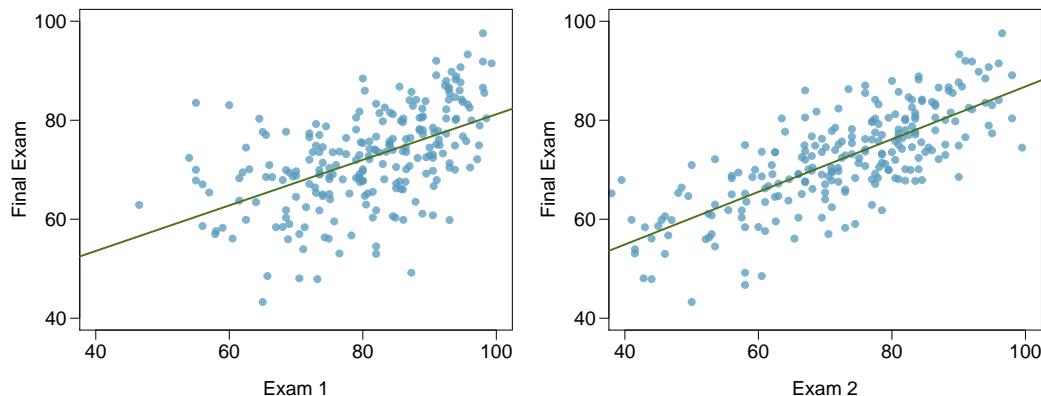


**7.4 Identify relationships, Part II.** For each of the six plots, identify the strength of the relationship (e.g. weak, moderate, or strong) in the data and whether fitting a linear model would be reasonable.

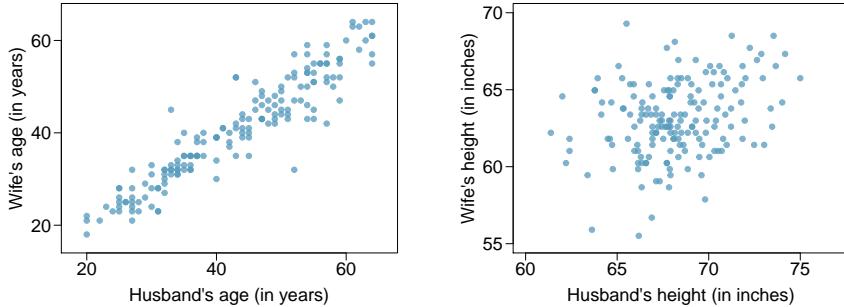


**7.5 Exams and grades.** The two scatterplots below show the relationship between final and mid-semester exam grades recorded during several years for a Statistics course at a university.

- Based on these graphs, which of the two exams has the strongest correlation with the final exam grade? Explain.
- Can you think of a reason why the correlation between the exam you chose in part (a) and the final exam is higher?



**7.6 Husbands and wives, Part I.** The Great Britain Office of Population Census and Surveys once collected data on a random sample of 170 married couples in Britain, recording the age (in years) and heights (converted here to inches) of the husbands and wives.<sup>16</sup> The scatterplot on the left shows the wife's age plotted against her husband's age, and the plot on the right shows wife's height plotted against husband's height.

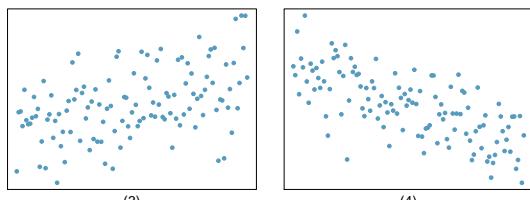
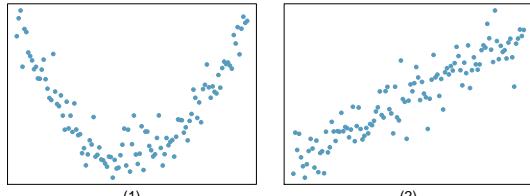


- (a) Describe the relationship between husbands' and wives' ages.
- (b) Describe the relationship between husbands' and wives' heights.
- (c) Which plot shows a stronger correlation? Explain your reasoning.
- (d) Data on heights were originally collected in centimeters, and then converted to inches. Does this conversion affect the correlation between husbands' and wives' heights?

**7.7 Match the correlation, Part I.**

Match the calculated correlations to the corresponding scatterplot.

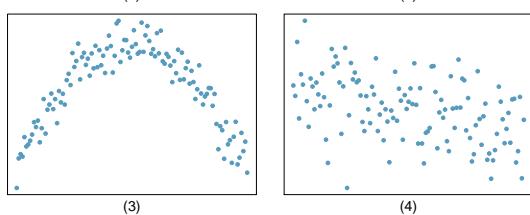
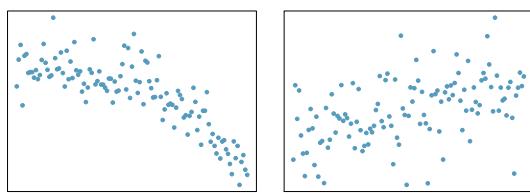
- (a)  $r = -0.7$
- (b)  $r = 0.45$
- (c)  $r = 0.06$
- (d)  $r = 0.92$



**7.8 Match the correlation, Part II.**

Match the calculated correlations to the corresponding scatterplot.

- (a)  $r = 0.49$
- (b)  $r = -0.48$
- (c)  $r = -0.03$
- (d)  $r = -0.85$



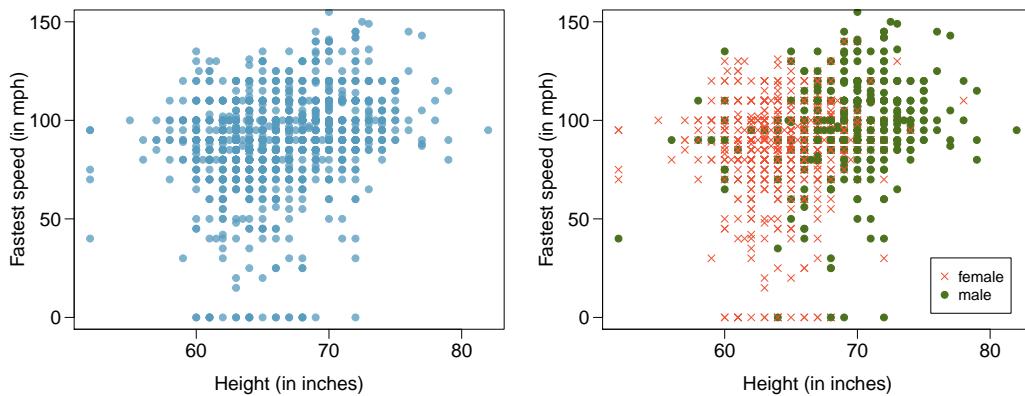
<sup>16</sup>D.J. Hand. *A handbook of small data sets*. Chapman & Hall/CRC, 1994.

**7.9 True / False.** Determine if the following statements are true or false. If false, explain why.

- A correlation coefficient of -0.90 indicates a stronger linear relationship than a correlation coefficient of 0.5.
- Correlation is a measure of the association between any two variables.

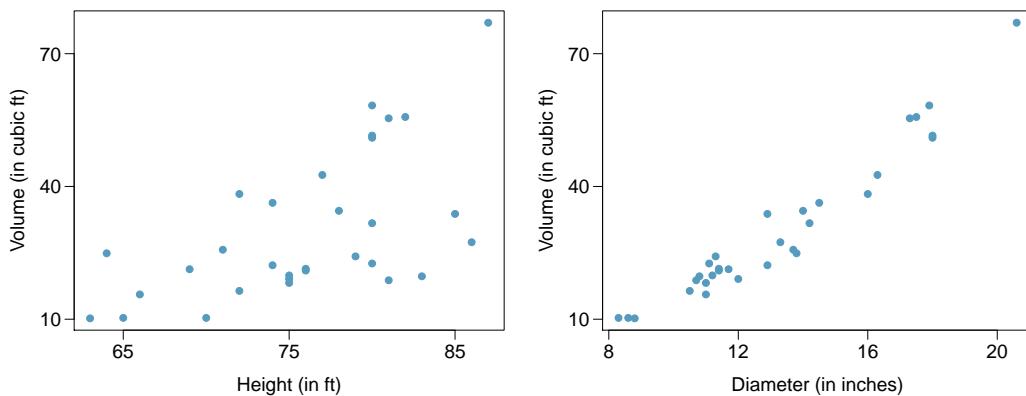
**7.10 Guess the correlation.** Eduardo and Rosie are both collecting data on number of rainy days in a year and the total rainfall for the year. Eduardo records rainfall in inches and Rosie in centimeters. How will their correlation coefficients compare?

**7.11 Speed and height.** 1,302 UCLA students were asked to fill out a survey where they were asked about their height, fastest speed they have ever driven, and gender. The scatterplot on the left displays the relationship between height and fastest speed, and the scatterplot on the right displays the breakdown by gender in this relationship.



- Describe the relationship between height and fastest speed.
- Why do you think these variables are positively associated?
- What role does gender play in the relationship between height and fastest driving speed?

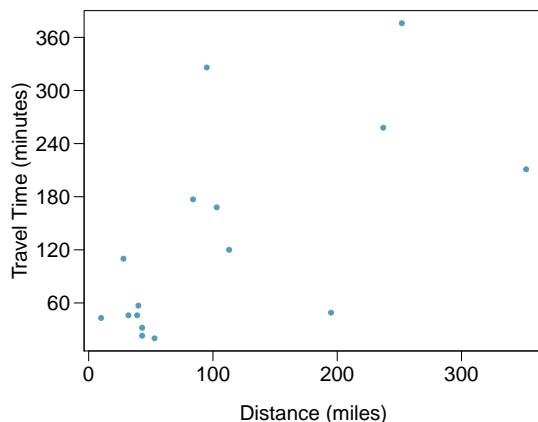
**7.12 Trees.** The scatterplots below show the relationship between height, diameter, and volume of timber in 31 felled black cherry trees. The diameter of the tree is measured 4.5 feet above the ground.<sup>17</sup>



- (a) Describe the relationship between volume and height of these trees.
- (b) Describe the relationship between volume and diameter of these trees.
- (c) Suppose you have height and diameter measurements for another black cherry tree. Which of these variables would be preferable to use to predict the volume of timber in this tree using a simple linear regression model? Explain your reasoning.

**7.13 The Coast Starlight, Part I.** The Coast Starlight Amtrak train runs from Seattle to Los Angeles. The scatterplot below displays the distance between each stop (in miles) and the amount of time it takes to travel from one stop to another (in minutes).

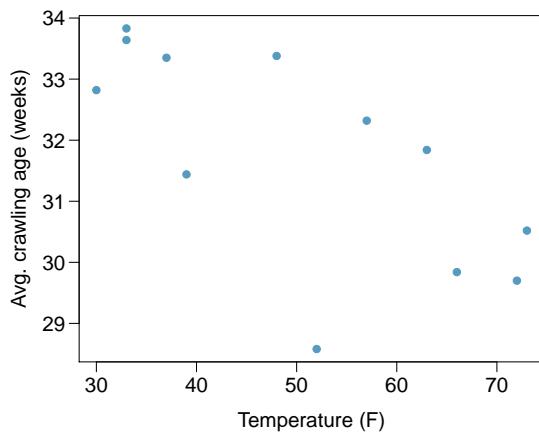
- (a) Describe the relationship between distance and travel time.
- (b) How would the relationship change if travel time was instead measured in hours, and distance was instead measured in kilometers?
- (c) Correlation between travel time (in miles) and distance (in minutes) is  $r = 0.636$ . What is the correlation between travel time (in kilometers) and distance (in hours)?



<sup>17</sup>Source: R Dataset, stat.ethz.ch/R-manual/R-patched/library/datasets/html/trees.html.

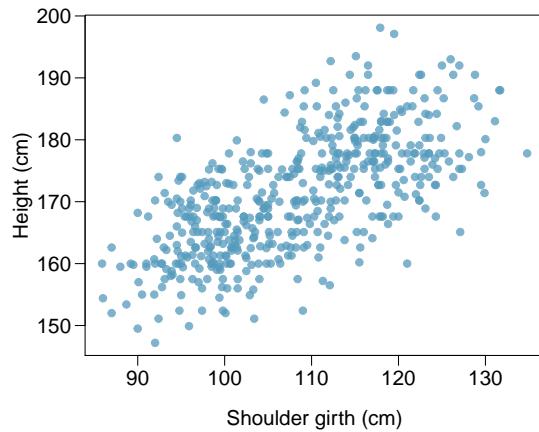
**7.14 Crawling babies, Part I.** A study conducted at the University of Denver investigated whether babies take longer to learn to crawl in cold months, when they are often bundled in clothes that restrict their movement, than in warmer months.<sup>18</sup> Infants born during the study year were split into twelve groups, one for each birth month. We consider the average crawling age of babies in each group against the average temperature when the babies are six months old (that's when babies often begin trying to crawl). Temperature is measured in degrees Fahrenheit ( $^{\circ}\text{F}$ ) and age is measured in weeks.

- (a) Describe the relationship between temperature and crawling age.
- (b) How would the relationship change if temperature was measured in degrees Celsius ( $^{\circ}\text{C}$ ) and age was measured in months?
- (c) The correlation between temperature in  $^{\circ}\text{F}$  and age in weeks was  $r = -0.70$ . If we converted the temperature to  $^{\circ}\text{C}$  and age to months, what would the correlation be?



**7.15 Body measurements, Part I.** Researchers studying anthropometry collected body girth measurements and skeletal diameter measurements, as well as age, weight, height and gender for 507 physically active individuals.<sup>19</sup> The scatterplot below shows the relationship between height and shoulder girth (over deltoid muscles), both measured in centimeters.

- (a) Describe the relationship between shoulder girth and height.
- (b) How would the relationship change if shoulder girth was measured in inches while the units of height remained in centimeters?

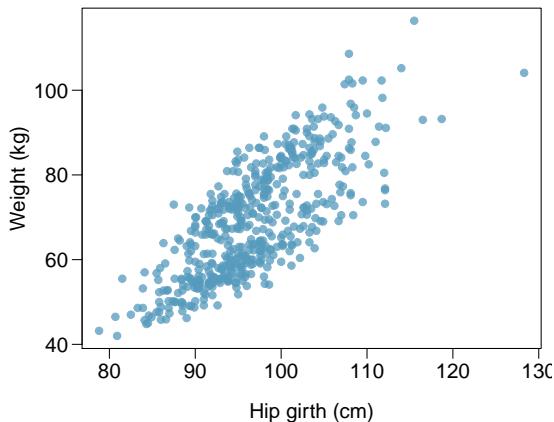


<sup>18</sup>J.B. Benson. "Season of birth and onset of locomotion: Theoretical and methodological implications". In: *Infant behavior and development* 16.1 (1993), pp. 69–81. ISSN: 0163-6383.

<sup>19</sup>G. Heinz et al. "Exploring relationships in body dimensions". In: *Journal of Statistics Education* 11.2 (2003).

**7.16 Body measurements, Part II.** The scatterplot below shows the relationship between weight measured in kilograms and hip girth measured in centimeters from the data described in Exercise 7.15.

- (a) Describe the relationship between hip girth and weight.
- (b) How would the relationship change if weight was measured in pounds while the units for hip girth remained in centimeters?



**7.17 Correlation, Part I.** What would be the correlation between the ages of husbands and wives if men always married women who were

- (a) 3 years younger than themselves?
- (b) 2 years older than themselves?
- (c) half as old as themselves?

**7.18 Correlation, Part II.** What would be the correlation between the annual salaries of males and females at a company if for a certain type of position men always made

- (a) \$5,000 more than women?
- (b) 25% more than women?
- (c) 15% less than women?

### 7.5.2 Fitting a line by least squares regression

**7.19 Units of regression.** Consider a regression predicting weight (kg) from height (cm) for a sample of adult males. What are the units of the correlation coefficient, the intercept, and the slope?

**7.20 Which is higher?** Determine if I or II is higher or if they are equal. Explain your reasoning.

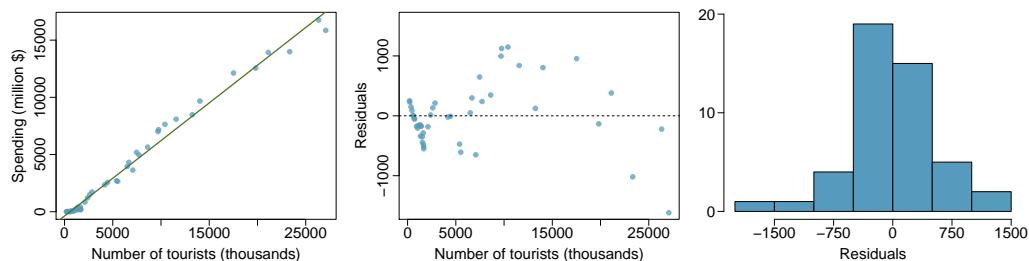
For a regression line, the uncertainty associated with the slope estimate,  $b_1$ , is higher when

- I. there is a lot of scatter around the regression line or
- II. there is very little scatter around the regression line

**7.21 Over-under, Part I.** Suppose we fit a regression line to predict the shelf life of an apple based on its weight. For a particular apple, we predict the shelf life to be 4.6 days. The apple's residual is -0.6 days. Did we over or under estimate the shelf-life of the apple? Explain your reasoning.

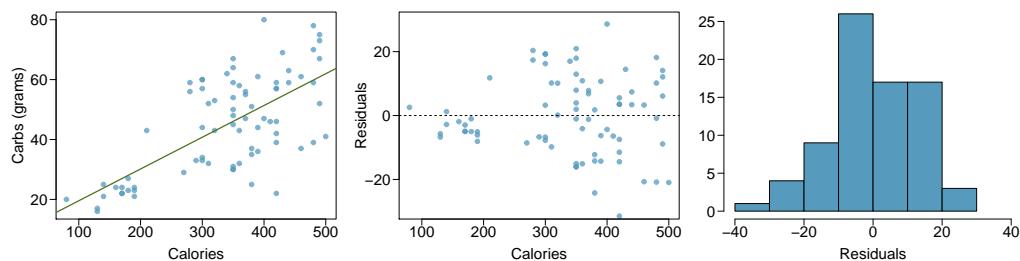
**7.22 Over-under, Part II.** Suppose we fit a regression line to predict the number of incidents of skin cancer per 1,000 people from the number of sunny days in a year. For a particular year, we predict the incidence of skin cancer to be 1.5 per 1,000 people, and the residual for this year is 0.5. Did we over or under estimate the incidence of skin cancer? Explain your reasoning.

**7.23 Tourism spending.** The Association of Turkish Travel Agencies reports the number of foreign tourists visiting Turkey and tourist spending by year.<sup>20</sup> Three plots are provided: scatterplot showing the relationship between these two variables along with the least squares fit, residuals plot, and histogram of residuals.



- Describe the relationship between number of tourists and spending.
- What are the explanatory and response variables?
- Why might we want to fit a regression line to these data?
- Do the data meet the conditions required for fitting a least squares line? In addition to the scatterplot, use the residual plot and histogram to answer this question.

**7.24 Nutrition at Starbucks, Part I.** The scatterplot below shows the relationship between the number of calories and amount of carbohydrates (in grams) Starbucks food menu items contain.<sup>21</sup> Since Starbucks only lists the number of calories on the display items, we are interested in predicting the amount of carbs a menu item has based on its calorie content.



- Describe the relationship between number of calories and amount of carbohydrates (in grams) that Starbucks food menu items contain.
- In this scenario, what are the explanatory and response variables?
- Why might we want to fit a regression line to these data?
- Do these data meet the conditions required for fitting a least squares line?

<sup>20</sup>Association of Turkish Travel Agencies, Foreign Visitors Figure & Tourist Spendings By Years.

<sup>21</sup>Source: Starbucks.com, collected on March 10, 2011,  
www.starbucks.com/menu/nutrition.



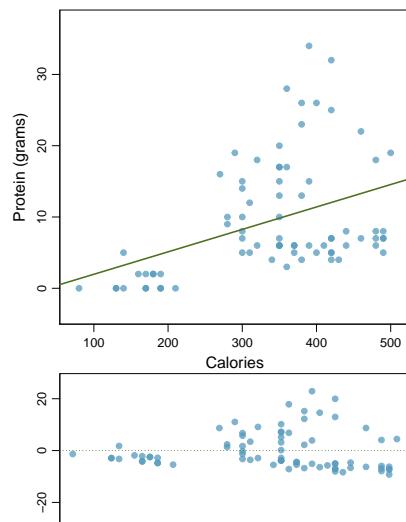
**7.25 The Coast Starlight, Part II.** Exercise 7.13 introduces data on the Coast Starlight Amtrak train that runs from Seattle to Los Angeles. The mean travel time from one stop to the next on the Coast Starlight is 129 mins, with a standard deviation of 113 minutes. The mean distance traveled from one stop to the next is 108 miles with a standard deviation of 99 miles. The correlation between travel time and distance is 0.636.

- Write the equation of the regression line for predicting travel time.
- Interpret the slope and the intercept in this context.
- Calculate  $R^2$  of the regression line for predicting travel time from distance traveled for the Coast Starlight, and interpret  $R^2$  in the context of the application.
- The distance between Santa Barbara and Los Angeles is 103 miles. Use the model to estimate the time it takes for the Starlight to travel between these two cities.
- It actually takes the Coast Starlight about 168 mins to travel from Santa Barbara to Los Angeles. Calculate the residual and explain the meaning of this residual value.
- Suppose Amtrak is considering adding a stop to the Coast Starlight 500 miles away from Los Angeles. Would it be appropriate to use this linear model to predict the travel time from Los Angeles to this point?

**7.26 Body measurements, Part III.** Exercise 7.15 introduces data on shoulder girth and height of a group of individuals. The mean shoulder girth is 107.20 cm with a standard deviation of 10.37 cm. The mean height is 171.14 cm with a standard deviation of 9.41 cm. The correlation between height and shoulder girth is 0.67.

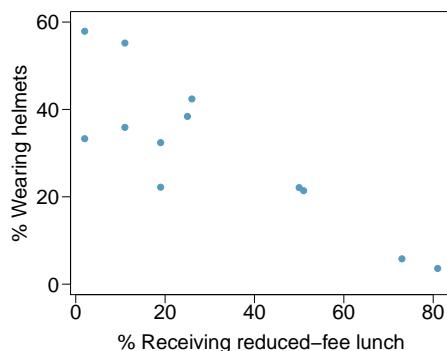
- Write the equation of the regression line for predicting height.
- Interpret the slope and the intercept in this context.
- Calculate  $R^2$  of the regression line for predicting height from shoulder girth, and interpret it in the context of the application.
- A randomly selected student from your class has a shoulder girth of 100 cm. Predict the height of this student using the model.
- The student from part (d) is 160 cm tall. Calculate the residual, and explain what this residual means.
- A one year old has a shoulder girth of 56 cm. Would it be appropriate to use this linear model to predict the height of this child?

**7.27 Nutrition at Starbucks, Part II.** Exercise 7.24 introduced a data set on nutrition information on Starbucks food menu items. Based on the scatterplot and the residual plot provided, describe the relationship between the protein content and calories of these menu items, and determine if a simple linear model is appropriate to predict amount of protein from the number of calories.



**7.28 Helmets and lunches.** The scatterplot shows the relationship between socioeconomic status measured as the percentage of children in a neighborhood receiving reduced-fee lunches at school (`lunch`) and the percentage of bike riders in the neighborhood wearing helmets (`helmet`). The average percentage of children receiving reduced-fee lunches is 30.8% with a standard deviation of 26.7% and the average percentage of bike riders wearing helmets is 38.8% with a standard deviation of 16.9%.

- If the  $R^2$  for the least-squares regression line for these data is 72%, what is the correlation between `lunch` and `helmet`?
- Calculate the slope and intercept for the least-squares regression line for these data.
- Interpret the intercept of the least-squares regression line in the context of the application.
- Interpret the slope of the least-squares regression line in the context of the application.
- What would the value of the residual be for a neighborhood where 40% of the children receive reduced-fee lunches and 40% of the bike riders wear helmets? Interpret the meaning of this residual in the context of the application.



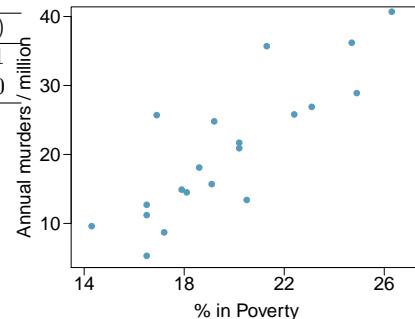
**7.29 Murders and poverty, Part I.** The following regression output is for predicting annual murders per million from percentage living in poverty in a random sample of 20 metropolitan areas.



	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-29.901	7.789	-3.839	0.001
poverty%	2.559	0.390	6.562	0.000

$s = 5.512 \quad R^2 = 70.52\% \quad R_{adj}^2 = 68.89\%$

- Write out the linear model.
- Interpret the intercept.
- Interpret the slope.
- Interpret  $R^2$ .
- Calculate the correlation coefficient.

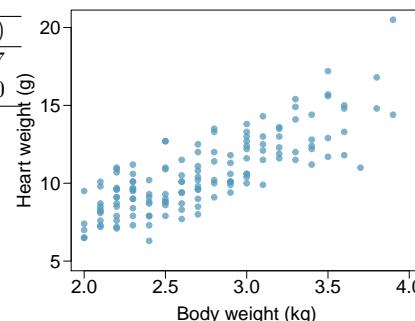


**7.30 Cats, Part I.** The following regression output is for predicting the heart weight (in g) of cats from their body weight (in kg). The coefficients are estimated using a dataset of 144 domestic cats.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-0.357	0.692	-0.515	0.607
body wt	4.034	0.250	16.119	0.000

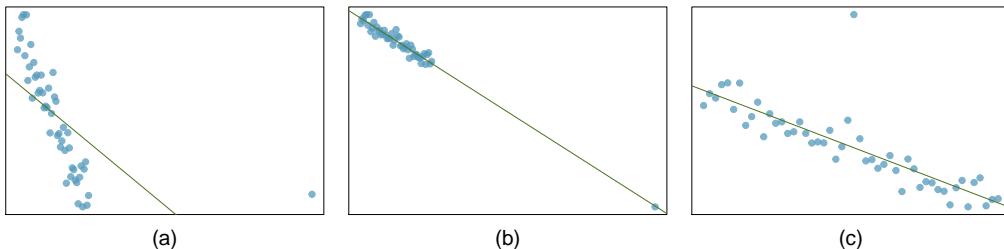
$s = 1.452 \quad R^2 = 64.66\% \quad R_{adj}^2 = 64.41\%$

- Write out the linear model.
- Interpret the intercept.
- Interpret the slope.
- Interpret  $R^2$ .
- Calculate the correlation coefficient.

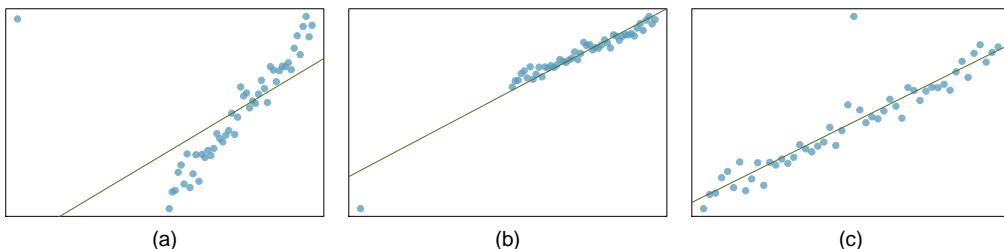


### 7.5.3 Types of outliers in linear regression

**7.31 Outliers, Part I.** Identify the outliers in the scatterplots shown below, and determine what type of outliers they are. Explain your reasoning.

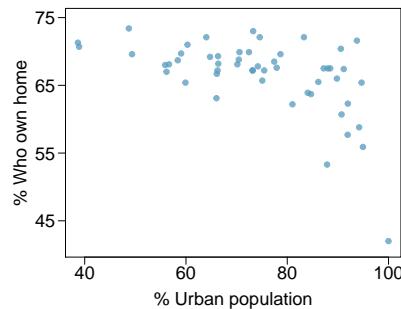


**7.32 Outliers, Part II.** Identify the outliers in the scatterplots shown below and determine what type of outliers they are. Explain your reasoning.



**7.33 Urban homeowners, Part I.** The scatterplot below shows the percent of families who own their home vs. the percent of the population living in urban areas in 2010.<sup>22</sup> There are 52 observations, each corresponding to a state in the US. Puerto Rico and District of Columbia are also included.

- Describe the relationship between the percent of families who own their home and the percent of the population living in urban areas in 2010.
- The outlier at the bottom right corner is District of Columbia, where 100% of the population is considered urban. What type of outlier is this observation?



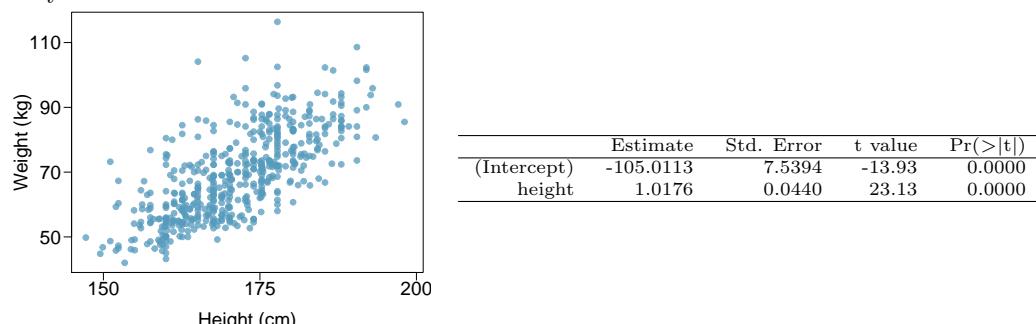
**7.34 Crawling babies, Part II.** Exercise 7.14 introduces data on the average monthly temperature during the month babies first try to crawl (about 6 months after birth) and the average first crawling age for babies born in a given month. A scatterplot of these two variables reveals a potential outlying month when the average temperature is about 53°F and average crawling age is about 28.5 weeks. Does this point have high leverage? Is it an influential point?

<sup>22</sup>United States Census Bureau, 2010 Census Urban and Rural Classification and Urban Area Criteria and Housing Characteristics: 2010.

### 7.5.4 Inference for linear regression

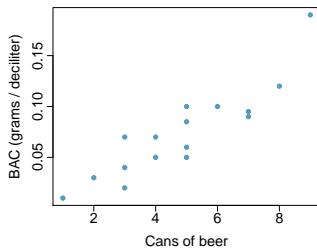
In the following exercises, visually check the conditions for fitting a least squares regression line, but you do not need to report these conditions in your solutions.

**7.35 Body measurements, Part IV.** The scatterplot and least squares summary below show the relationship between weight measured in kilograms and height measured in centimeters of 507 physically active individuals.



- Describe the relationship between height and weight.
- Write the equation of the regression line. Interpret the slope and intercept in context.
- Do the data provide strong evidence that an increase in height is associated with an increase in weight? State the null and alternative hypotheses, report the p-value, and state your conclusion.
- The correlation coefficient for height and weight is 0.72. Calculate  $R^2$  and interpret it in context.

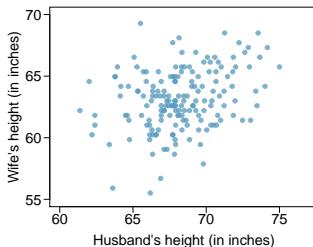
**7.36 Beer and blood alcohol content.** Many people believe that gender, weight, drinking habits, and many other factors are much more important in predicting blood alcohol content (BAC) than simply considering the number of drinks a person consumed. Here we examine data from sixteen student volunteers at Ohio State University who each drank a randomly assigned number of cans of beer. These students were evenly divided between men and women, and they differed in weight and drinking habits. Thirty minutes later, a police officer measured their blood alcohol content (BAC) in grams of alcohol per deciliter of blood.<sup>23</sup> The scatterplot and regression table summarize the findings.



	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-0.0127	0.0126	-1.00	0.3320
beers	0.0180	0.0024	7.48	0.0000

- Describe the relationship between the number of cans of beer and BAC.
- Write the equation of the regression line. Interpret the slope and intercept in context.
- Do the data provide strong evidence that drinking more cans of beer is associated with an increase in blood alcohol? State the null and alternative hypotheses, report the p-value, and state your conclusion.
- The correlation coefficient for number of cans of beer and BAC is 0.89. Calculate  $R^2$  and interpret it in context.
- Suppose we visit a bar, ask people how many drinks they have had, and also take their BAC. Do you think the relationship between number of drinks and BAC would be as strong as the relationship found in the Ohio State study?

**7.37 Husbands and wives, Part II.** The scatterplot below summarizes husbands' and wives' heights in a random sample of 170 married couples in Britain, where both partners' ages are below 65 years. Summary output of the least squares fit for predicting wife's height from husband's height is also provided in the table.

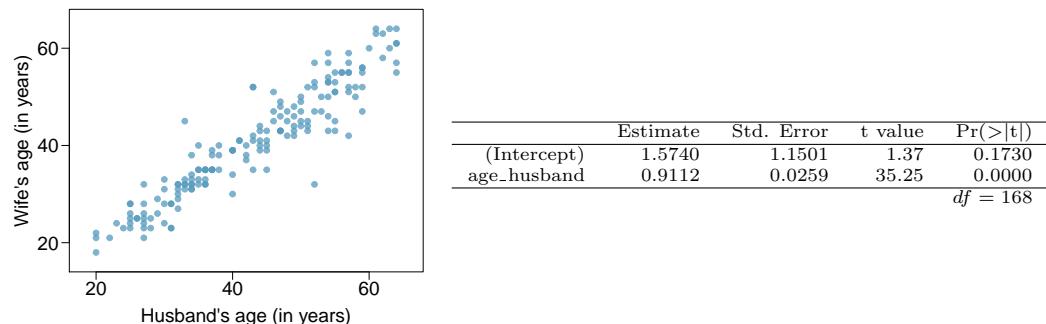


	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	43.5755	4.6842	9.30	0.0000
height_husband	0.2863	0.0686	4.17	0.0000

- Is there strong evidence that taller men marry taller women? State the hypotheses and include any information used to conduct the test.
- Write the equation of the regression line for predicting wife's height from husband's height.
- Interpret the slope and intercept in the context of the application.
- Given that  $R^2 = 0.09$ , what is the correlation of heights in this data set?
- You meet a married man from Britain who is 5'9" (69 inches). What would you predict his wife's height to be? How reliable is this prediction?
- You meet another married man from Britain who is 6'7" (79 inches). Would it be wise to use the same linear model to predict his wife's height? Why or why not?

<sup>23</sup>J. Malkevitch and L.M. Lesser. *For All Practical Purposes: Mathematical Literacy in Today's World*. WH Freeman & Co, 2008.

**7.38 Husbands and wives, Part III.** Exercise 7.37 presents a scatterplot displaying the relationship between husbands' and wives' ages in a random sample of 170 married couples in Britain, where both partners' ages are below 65 years. Given below is summary output of the least squares fit for predicting wife's age from husband's age.

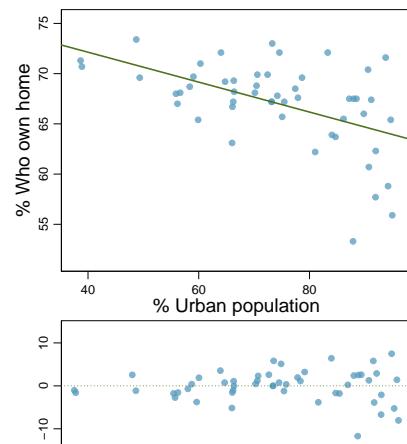


- We might wonder, is the age difference between husbands and wives consistent across ages? If this were the case, then the slope parameter would be  $\beta_1 = 1$ . Use the information above to evaluate if there is strong evidence that the difference in husband and wife ages differs for different ages.
- Write the equation of the regression line for predicting wife's age from husband's age.
- Interpret the slope and intercept in context.
- Given that  $R^2 = 0.88$ , what is the correlation of ages in this data set?
- You meet a married man from Britain who is 55 years old. What would you predict his wife's age to be? How reliable is this prediction?
- You meet another married man from Britain who is 85 years old. Would it be wise to use the same linear model to predict his wife's age? Explain.

### 7.39 Urban homeowners, Part II.

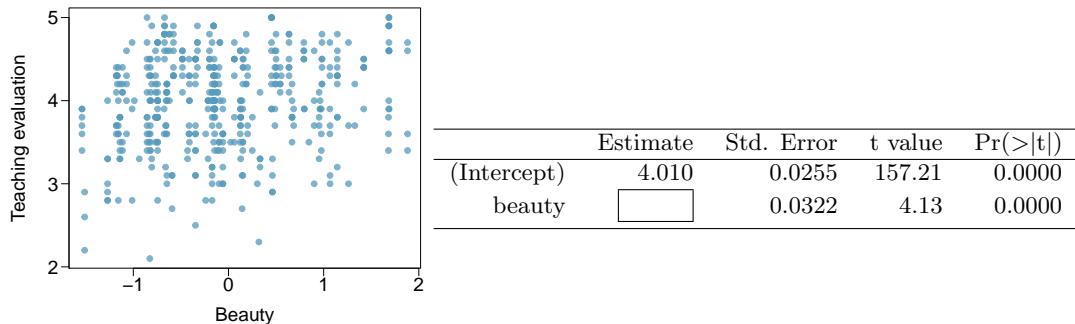
Exercise 7.33 gives a scatterplot displaying the relationship between the percent of families that own their home and the percent of the population living in urban areas. Below is a similar scatterplot, excluding District of Columbia, as well as the residuals plot. There were 51 cases.

- For these data,  $R^2 = 0.28$ . What is the correlation? How can you tell if it is positive or negative?
- Examine the residual plot. What do you observe? Is a simple least squares fit appropriate for these data?

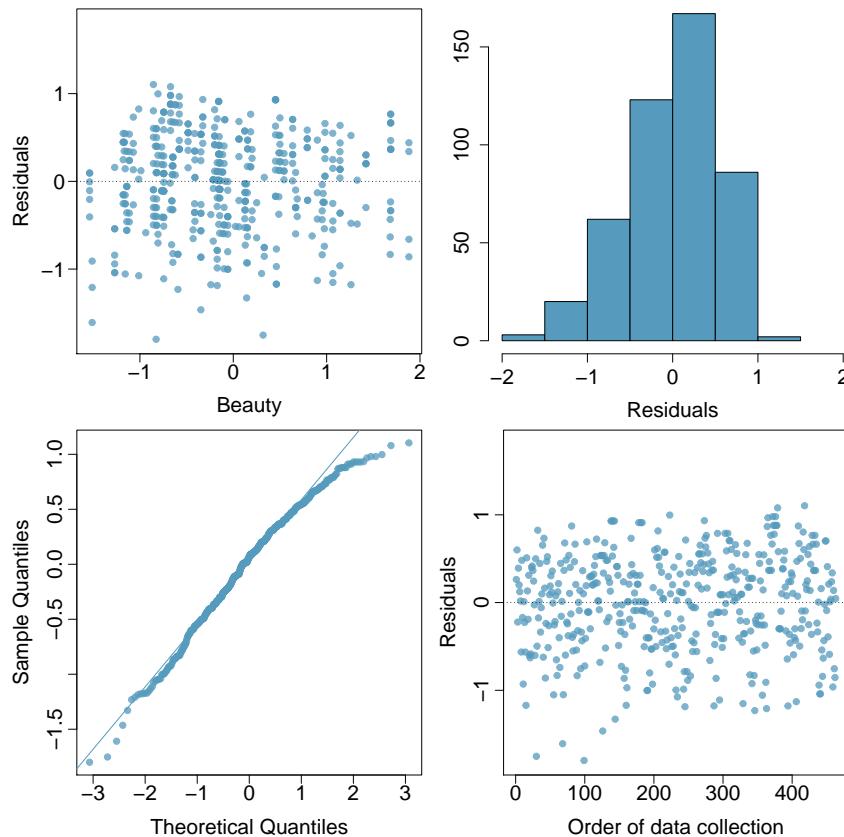


**7.40 Rate my professor.** Many college courses conclude by giving students the opportunity to evaluate the course and the instructor anonymously. However, the use of these student evaluations as an indicator of course quality and teaching effectiveness is often criticized because these measures may reflect the influence of non-teaching related characteristics, such as the physical appearance of the instructor. Researchers at University of Texas, Austin collected data on teaching evaluation score (higher score means better) and standardized beauty score (a score of 0 means average, negative score means below average, and a positive score means above average) for a

sample of 463 professors.<sup>24</sup> The scatterplot below shows the relationship between these variables, and also provided is a regression output for predicting teaching evaluation score from beauty score.



- Given that the average standardized beauty score is -0.0883 and average teaching evaluation score is 3.9983, calculate the slope. Alternatively, the slope may be computed using just the information provided in the model summary table.
- Do these data provide convincing evidence that the slope of the relationship between teaching evaluation and beauty is positive? Explain your reasoning.
- List the conditions required for linear regression and check if each one is satisfied for this model based on the following diagnostic plots.



<sup>24</sup>Daniel S Hamermesh and Amy Parker. "Beauty in the classroom: Instructors pulchritude and putative pedagogical productivity". In: *Economics of Education Review* 24.4 (2005), pp. 369–376.



**7.41 Murders and poverty, Part II.** Exercise 7.29 presents regression output from a model for predicting annual murders per million from percentage living in poverty based on a random sample of 20 metropolitan areas. The model output is also provided below.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-29.901	7.789	-3.839	0.001
poverty%	2.559	0.390	6.562	0.000
$s = 5.512$		$R^2 = 70.52\%$		$R^2_{adj} = 68.89\%$

- (a) What are the hypotheses for evaluating whether poverty percentage is a significant predictor of murder rate?
- (b) State the conclusion of the hypothesis test from part (a) in context of the data.
- (c) Calculate a 95% confidence interval for the slope of poverty percentage, and interpret it in context of the data.
- (d) Do your results from the hypothesis test and the confidence interval agree? Explain.

**7.42 Babies.** Is the gestational age (time between conception and birth) of a low birth-weight baby useful in predicting head circumference at birth? Twenty-five low birth-weight babies were studied at a Harvard teaching hospital; the investigators calculated the regression of head circumference (measured in centimeters) against gestational age (measured in weeks). The estimated regression line is

$$\widehat{\text{head\_circumference}} = 3.91 + 0.78 \times \text{gestational\_age}$$

- (a) What is the predicted head circumference for a baby whose gestational age is 28 weeks?
- (b) The standard error for the coefficient of gestational age is 0.35, which is associated with  $df = 23$ . Does the model provide strong evidence that gestational age is significantly associated with head circumference?

**7.43 Murders and poverty, Part III.** In Exercises 7.41 you evaluated whether poverty percentage is a significant predictor of murder rate. How, if at all, would your answer change if we wanted to find out whether poverty percentage is positively associated with murder rate. Make sure to include the appropriate p-value for this hypothesis test in your answer.

**7.44 Cats, Part II.** Exercise 7.30 presents regression output from a model for predicting the heart weight (in g) of cats from their body weight (in kg). The coefficients are estimated using a dataset of 144 domestic cats. The model output is also provided below.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-0.357	0.692	-0.515	0.607
body wt	4.034	0.250	16.119	0.000
$s = 1.452$		$R^2 = 64.66\%$		$R^2_{adj} = 64.41\%$

- (a) What are the hypotheses for evaluating whether body weight is positively associated with heart weight in cats?
- (b) State the conclusion of the hypothesis test from part (a) in context of the data.
- (c) Calculate a 95% confidence interval for the slope of body weight, and interpret it in context of the data.
- (d) Do your results from the hypothesis test and the confidence interval agree? Explain.

## Appendix A

# End of chapter exercise solutions

### 1 Introduction to data

**1.1** (a) Treatment:  $10/43 = 0.23 \rightarrow 23\%$ . Control:  $2/46 = 0.04 \rightarrow 4\%$ . (b) There is a 19% difference between the pain reduction rates in the two groups. At first glance, it appears patients in the treatment group are more likely to experience pain reduction from the acupuncture treatment. (c) Answers may vary but should be sensible. Two possible answers: <sup>1</sup>Though the groups' difference is big, I'm skeptical the results show a real difference and think this might be due to chance. <sup>2</sup>The difference in these rates looks pretty big, so I suspect acupuncture is having a positive impact on pain.

**1.3** (a) 143,196 eligible study subjects born in Southern California between 1989 and 1993. (b) Measurements of carbon monoxide, nitrogen dioxide, ozone, and particulate matter less than  $10\mu\text{g}/\text{m}^3$  ( $\text{PM}_{10}$ ) collected at air-quality-monitoring stations as well as length of gestation. Continuous numerical variables. (c) "Is there an association between air pollution exposure and preterm births?"

**1.5** (a) 160 children. (b) Age (numerical, continuous), sex (categorical), whether they were an only child or not (categorical), and whether they cheated or not (categorical). (c) Research question: "Does explicitly telling children not to cheat affect their likelihood to cheat?"

**1.7** (a)  $50 \times 3 = 150$ . (b) Four continuous numerical variables: sepal length, sepal width, petal length, and petal width. (c) One categorical variable, species, with three levels: *setosa*, *versicolor*, and *virginica*.

**1.9** (a) Population: all births, sample: 143,196 births between 1989 and 1993 in Southern California. (b) If births in this time span at the geography can be considered to be representative of all births, then the results are generalizable to the population of Southern California. However, since the study is observational the findings cannot be used to establish causal relationships.

**1.11** (a) Population: all asthma patients aged 18-69 who rely on medication for asthma treatment. Sample: 600 such patients. (b) If the patients in this sample, who are likely not randomly sampled, can be considered to be representative of all asthma patients aged 18-69 who rely on medication for asthma treatment, then the results are generalizable to the population defined above. Additionally, since the study is experimental, the findings can be used to establish causal relationships.

**1.13** (a) Observation. (b) Variable. (c) Sample statistic (mean). (d) Population parameter (mean).

**1.15** (a) Explanatory: number of study hours per week. Response: GPA. (b) Somewhat weak positive relationship with data becoming more sparse as the number of study hours increases. One responded reported a GPA above 4.0, which is clearly a data error. There are a few respondents who reported unusually high study hours (60 and 70 hours/week). Variability in GPA is much higher for students who study less than those who study more, which might be due to the fact that there aren't many respondents who reported studying higher hours. (c) Observational. (d) Since observational, cannot infer causation.

**1.17** (a) Observational. (b) Use stratified sampling to randomly sample a fixed number of students, say 10, from each section for a total sample size of 40 students.

**1.19** (a) Positive, non-linear, somewhat strong. Countries in which a higher percentage of the population have access to the internet also tend to have higher average life expectancies, however rise in life expectancy trails off before around 80 years old. (b) Observational. (c) Wealth: countries with individuals who can widely afford the internet can probably also afford basic medical care. (Note: Answers may vary.)

**1.21** (a) Simple random sampling is okay. In fact, it's rare for simple random sampling to not be a reasonable sampling method! (b) The student opinions may vary by field of study, so the stratifying by this variable makes sense and would be reasonable. (c) Students of similar ages are probably going to have more similar opinions, and we want clusters to be diverse with respect to the outcome of interest, so this would **not** be a good approach. (Additional thought: the clusters in this case may also have very different numbers of people, which can also create unexpected sample sizes.)

**1.23** (a) The cases are 200 randomly sampled men and women. (b) The response variable is attitude towards a fictional microwave oven. (c) The explanatory variable is dispositional attitude. (d) Yes, the cases are sampled randomly. (e) This is an observational study since there is no random assignment to treatments. (f) No, we cannot establish a causal link between the ex-

planatory and response variables since the study is observational. (g) Yes, the results of the study can be generalized to the population at large since the sample is random.

**1.25** (a) Non-responders may have a different response to this question, e.g. parents who returned the surveys likely don't have difficulty spending time with their children. (b) It is unlikely that the women who were reached at the same address 3 years later are a random sample. These missing responders are probably renters (as opposed to homeowners) which means that they might be in a lower socio-economic status than the respondents. (c) There is no control group in this study, this is an observational study, and there may be confounding variables, e.g. these people may go running because they are generally healthier and/or do other exercises.

**1.27** (a) Simple random sample. Non-response bias, if only those people who have strong opinions about the survey responds his sample may not be representative of the population. (b) Convenience sample. Under coverage bias, his sample may not be representative of the population since it consists only of his friends. It is also possible that the study will have non-response bias if some choose to not bring back the survey. (c) Convenience sample. This will have a similar issues to handing out surveys to friends. (d) Multi-stage sampling. If the classes are similar to each other with respect to student composition this approach should not introduce bias, other than potential non-response bias.

**1.29** No, students were not randomly sampled (voluntary sample) and the sample only contains college students at a university in Ontario.

**1.31** (a) Exam performance. (b) Light level: fluorescent overhead lighting, yellow overhead lighting, no overhead lighting (only desk lamps). (c) Sex: man, woman.

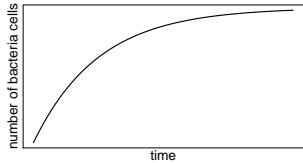
**1.33** (a) Exam performance. (b) Light level (overhead lighting, yellow overhead lighting, no overhead lighting) and noise level (no noise, construction noise, and human chatter noise). (c) Since the researchers want to ensure equal gender representation, sex will be a blocking variable.

**1.35** Need randomization and blinding. One possible outline: (1) Prepare two cups for each participant, one containing regular Coke and the other containing Diet Coke. Make sure the cups are identical and contain equal amounts of soda. Label the cups A (regular) and B (diet). (Be sure to randomize A and B for each trial!) (2) Give each participant the two cups, one cup at a time, in random order, and ask the participant to record a value that indicates how much she liked the beverage. Be sure that neither the participant nor the person handing out the cups knows the identity of the beverage to make this a double-blind experiment. (Answers may vary.)

**1.37** (a) Experiment. (b) Treatment: 25 grams of chia seeds twice a day, control: placebo. (c) Yes, gender. (d) Yes, single blind since the patients were blinded to the treatment they received. (e) Since this is an experiment, we can make a causal statement. However, since the sample is not random, the causal statement cannot be generalized to the population at large.

**1.39** (a) 1: linear. 3: nonlinear.  
(b) 4: linear. (c) 2.

**1.41**



**1.43** (a) Population mean,  $\mu_{2007} = 52$ ; sample mean,  $\bar{x}_{2008} = 58$ . (b) Population mean,  $\mu_{2001} = 3.37$ ; sample mean,  $\bar{x}_{2012} = 3.59$ .

**1.45** Any 10 employees whose average number of days off is between the minimum and the mean number of days off for the entire workforce at this plant.

**1.47** (a) Dist 2 has a higher mean since  $20 > 13$ , and a higher standard deviation since 20 is further from the rest of the data than 13. (b) Dist 1 has a higher mean since  $-20 > -40$ , and Dist 2 has a higher standard deviation since -40 is farther away from the rest of the data than -20. (c) Dist 2 has a higher mean since all values in this distribution are higher than those in Dist 1, but both distribution have the same standard deviation since they are equally variable around their respective means. (d) Both

distributions have the same mean since they're both centered at 300, but Dist 2 has a higher standard deviation since the observations are farther from the mean than in Dist 1.

**1.49** (a)  $Q1 \approx 5$ , median  $\approx 15$ ,  $Q3 \approx 35$   
(b) Since the distribution is right skewed, we would expect the mean to be higher than the median.

**1.51** (a) About 30. (b) Since the distribution is right skewed the mean is higher than the median. (c) Q1: between 15 and 20, Q3: between 35 and 40, IQR: about 20. (d) Values that are considered to be unusually low or high lie more than  $1.5 \times \text{IQR}$  away from the quartiles. Upper fence:  $Q3 + 1.5 \times \text{IQR} = 37.5 + 1.5 \times 20 = 67.5$ ; Lower fence:  $Q1 - 1.5 \times \text{IQR} = 17.5 + 1.5 \times 20 = -12.5$ ; The lowest AQI recorded is not lower than 5 and the highest AQI recorded is not higher than 65, which are both within the fences. Therefore none of the days in this sample would be considered to have an unusually low or high AQI.

**1.53** The histogram shows that the distribution is bimodal, which is not apparent in the box plot. The box plot makes it easy to identify more precise values of observations outside of the whiskers.

**1.55** (a) The distribution of number of pets per household is likely right skewed as there is a natural boundary at 0 and only a few people have many pets. Therefore the center would be best described by the median, and variability would be best described by the IQR. (b) The distribution of number of distance to work is likely right skewed as there is a natural boundary at 0 and only a few people live a very long distance from work. Therefore the center would be best described by the median, and variability would be best described by the IQR. (c) The distribution of heights of males is likely symmetric. Therefore the center would be best described by the mean, and variability would be best described by the standard deviation.

**1.57** No, we would expect this distribution to be right skewed. There are two reasons for this: (1) there is a natural boundary at 0 (it is not possible to watch less than 0 hours of TV), (2) the standard deviation of the distribution is very large compared to the mean.

**1.59** The statement “50% of Facebook users have over 100 friends” means that the median number of friends is 100, which is lower than the mean number of friends (190), which suggests a right skewed distribution for the number of friends of Facebook users.

**1.61** (a) The median is a much better measure of the typical amount earned by these 42 people. The mean is much higher than the income of 40 of the 42 people. This is because the mean is an arithmetic average and gets affected by the two extreme observations. The median does not get effected as much since it is robust to outliers. (b) The IQR is a much better measure of variability in the amounts earned by nearly all of the 42 people. The standard deviation gets affected greatly by the two high salaries, but the IQR is robust to these extreme observations.

**1.63** (a) The distribution is unimodal and symmetric with a mean of about 25 minutes and a standard deviation of about 5 minutes. There does not appear to be any counties with unusually high or low mean travel times. Since the distribution is already unimodal and symmetric, a log transformation is not necessary. (b) Answers will vary. There are pockets of longer travel time around DC, Southeastern NY, Chicago, Minneapolis, Los Angeles, and many other big cities. There is also a large section of shorter average commute times that overlap with farmland in the Midwest. Many farmers' homes are adjacent to their farmland, so their commute would be brief, which may explain why the average commute time for these counties is relatively low.

**1.65** (a) We see the order of the categories and the relative frequencies in the bar plot. (b) There are no features that are apparent in the pie chart but not in the bar plot. (c) We usually prefer to use a bar plot as we can also see the relative frequencies of the categories in this graph.

**1.67** The vertical locations at which the ideological groups break into the Yes, No, and Not Sure categories differ, which indicates that like-

lihood of supporting the DREAM act varies by political ideology. This suggests that the two variables may be dependent.

**1.69** (a) (i) False. Instead of comparing counts, we should compare percentages of people in each group who suffered cardiovascular problems. (ii) True. (iii) False. Association does not imply causation. We cannot infer a causal relationship based on an observational study. The difference from part (ii) is subtle. (iv) True.

(b) Proportion of all patients who had cardiovascular problems:  $\frac{7,979}{227,571} \approx 0.035$

(c) The expected number of heart attacks in the rosiglitazone group, if having cardiovascular problems and treatment were independent, can be calculated as the number of patients in that group multiplied by the overall cardiovascular problem rate in the study:  $67,593 * \frac{7,979}{227,571} \approx 2370$ .

(d) (i)  $H_0$ : The treatment and cardiovascular problems are independent. They have no relationship, and the difference in incidence rates between the rosiglitazone and pioglitazone groups is due to chance.  $H_A$ : The treatment and cardiovascular problems are not independent. The difference in the incidence rates between the rosiglitazone and pioglitazone groups is not due to chance and rosiglitazone is associated with an increased risk of serious cardiovascular problems. (ii) A higher number of patients with cardiovascular problems than expected under the assumption of independence would provide support for the alternative hypothesis as this would suggest that rosiglitazone increases the risk of such problems. (iii) In the actual study, we observed 2,593 cardiovascular events in the rosiglitazone group. In the 1,000 simulations under the independence model, we observed somewhat less than 2,593 in every single simulation, which suggests that the actual results did not come from the independence model. That is, the variables do not appear to be independent, and we reject the independence model in favor of the alternative. The study's results provide convincing evidence that rosiglitazone is associated with an increased risk of cardiovascular problems.

## 2 Probability

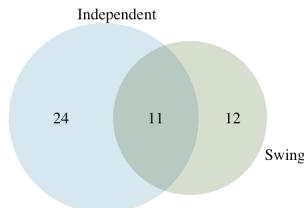
**??** (a) False. These are independent trials. (b) False. There are red face cards. (c) True. A card cannot be both a face card and an ace.

**??** (a) 10 tosses. Fewer tosses mean more variability in the sample fraction of heads, meaning there's a better chance of getting at least 60% heads. (b) 100 tosses. More flips means the observed proportion of heads would often be closer to the average, 0.50, and therefore also above 0.40. (c) 100 tosses. With more flips, the observed proportion of heads would often be closer to the average, 0.50. (d) 10 tosses. Fewer flips would increase variability in the fraction of tosses that are heads.

**??** (a)  $0.5^{10} = 0.00098$ . (b)  $0.5^{10} = 0.00098$ . (c)  $P(\text{at least one tails}) = 1 - P(\text{no tails}) = 1 - (0.5^{10}) \approx 1 - 0.001 = 0.999$ .

**??** (a) No, there are voters who are both independent and swing voters.

(b)



(c) Each Independent voter is either a swing voter or not. Since 35% of voters are Independents and 11% are both Independent and swing voters, the other 24% must not be swing voters. (d) 0.47. (e) 0.53. (f)  $P(\text{Independent}) \times P(\text{swing}) = 0.35 \times 0.23 = 0.08$ , which does not equal  $P(\text{Independent and swing}) = 0.11$ , so the events are dependent.

**??** (a) If the class is not graded on a curve, they are independent. If graded on a curve, then neither independent nor disjoint – unless the instructor will only give one A, which is a situation we will ignore in parts (b) and (c). (b) They are probably not independent: if you study together, your study habits would be related, which suggests your course performances are also related. (c) No. See the answer to part (a) when the course is not graded on a curve. More generally: if two things are un-

related (independent), then one occurring does not preclude the other from occurring.

**??** (a)  $0.16 + 0.09 = 0.25$ . (b)  $0.17 + 0.09 = 0.26$ . (c) Assuming that the education level of the husband and wife are independent:  $0.25 \times 0.26 = 0.065$ . You might also notice we actually made a second assumption: that the decision to get married is unrelated to education level. (d) The husband/wife independence assumption is probably not reasonable, because people often marry another person with a comparable level of education. We will leave it to you to think about whether the second assumption noted in part (c) is reasonable.

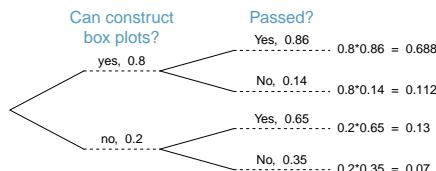
**??** (a) Invalid. Sum is greater than 1. (b) Valid. Probabilities are between 0 and 1, and they sum to 1. In this class, every student gets a C. (c) Invalid. Sum is less than 1. (d) Invalid. There is a negative probability. (e) Valid. Probabilities are between 0 and 1, and they sum to 1. (f) Invalid. There is a negative probability.

**??** (a) No, but we could if A and B are independent. (b-i) 0.21. (b-ii) 0.79. (b-iii) 0.3. (c) No, because  $0.1 \neq 0.21$ , where 0.21 was the value computed under independence from part (a). (d) 0.143.

**??** (a) No, 0.18 of respondents fall into this combination. (b)  $0.60 + 0.20 - 0.18 = 0.62$ . (c)  $0.18/0.20 = 0.9$ . (d)  $0.11/0.33 \approx 0.33$ . (e) No, otherwise the answers to (c) and (d) would be the same. (f)  $0.06/0.34 \approx 0.18$ .

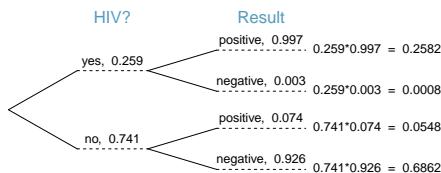
**??** (a) No. There are 6 females who like Five Guys Burgers. (b)  $162/248 = 0.65$ . (c)  $181/252 = 0.72$ . (d) Under the assumption of a dating choices being independent of hamburger preference, which on the surface seems reasonable:  $0.65 \times 0.72 = 0.468$ . (e)  $(252 + 6 - 1)/500 = 0.514$ .

**??** (a)

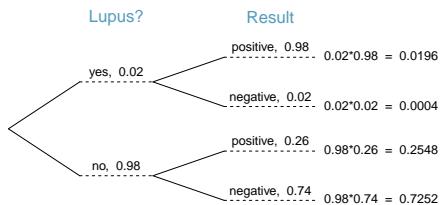


(b) 0.84

?? 0.8247.



?? 0.0714. Even when a patient tests positive for lupus, there is only a 7.14% chance that he actually has lupus. House may be right.

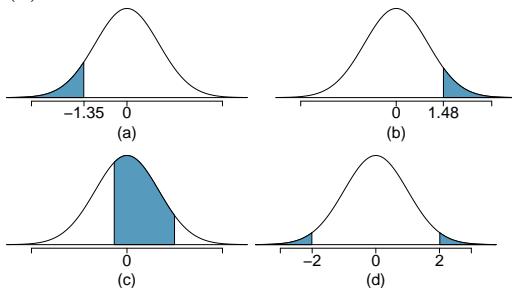


?? (a) 0.3. (b) 0.3. (c) 0.3. (d)  $0.3 \times 0.3 = 0.09$ .  
 (e) Yes, the population that is being sampled from is identical in each draw.

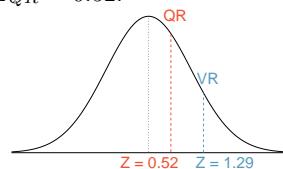
?? (a)  $2/9 \approx 0.22$ . (b)  $3/9 \approx 0.33$ . (c)  $\frac{3}{10} \times \frac{2}{9} \approx 0.067$ . (d) No, e.g. in this exercise, removing one chip meaningfully changes the probability

### 3 Distributions of random variables

3.1 (a) 8.85%. (b) 6.94%. (c) 58.86%.  
 (d) 4.56%.



3.3 (a) Verbal:  $N(\mu = 151, \sigma = 7)$ , Quant:  $N(\mu = 153, \sigma = 7.67)$ . (b)  $Z_{VR} = 1.29$ ,  $Z_{QR} = 0.52$ .



(c) She scored 1.29 standard deviations above

of what might be drawn next.

??  $P(1^{\text{st}} \text{leggings}, 2^{\text{nd}} \text{jeans}, 3^{\text{rd}} \text{jeans}) = \frac{5}{24} \times \frac{7}{23} \times \frac{6}{22} = 0.0173$ . However, the person with leggings could have come 2nd or 3rd, and these each have this same probability, so  $3 \times 0.0173 = 0.0519$ .

?? (a) 13. (b) No, these 27 students are not a random sample from the university's student population. For example, it might be argued that the proportion of smokers among students who go to the gym at 9 am on a Saturday morning would be lower than the proportion of smokers in the university as a whole.

?? (a)  $E(X) = 3.59$ .  $SD(X) = 9.64$ . (b)  $E(X) = -1.41$ .  $SD(X) = 9.64$ . (c) No, the expected net profit is negative, so on average you expect to lose money.

?? 5% increase in value.

??  $E = -0.0526$ .  $SD = 0.9986$ .

?? (a)  $E = \$3.90$ .  $SD = \$0.34$ .  
 (b)  $E = \$27.30$ .  $SD = \$0.89$ .

?? Approximate answers are OK.

(a)  $(29 + 32)/144 = 0.42$ . (b)  $21/144 = 0.15$ .  
 (c)  $(26 + 12 + 15)/144 = 0.37$ .

the mean on the Verbal Reasoning section and 0.52 standard deviations above the mean on the Quantitative Reasoning section. (d) She did better on the Verbal Reasoning section since her Z-score on that section was higher. (e)  $Perc_{VR} = 0.9007 \approx 90\%$ ,  $Perc_{QR} = 0.6990 \approx 70\%$ . (f)  $100\% - 90\% = 10\%$  did better than her on VR, and  $100\% - 70\% = 30\%$  did better than her on QR. (g) We cannot compare the raw scores since they are on different scales. Comparing her percentile scores is more appropriate when comparing her performance to others. (h) Answer to part (b) would not change as Z-scores can be calculated for distributions that are not normal. However, we could not answer parts (d)-(f) since we cannot use the normal probability table to calculate probabilities and percentiles without a normal model.

3.5 (a)  $Z = 0.84$ , which corresponds to approximately 160 on QR. (b)  $Z = -0.52$ , which corresponds to approximately 147 on VR.

3.7 (a)  $Z = 1.2 \rightarrow 0.1151$ .

(b)  $Z = -1.28 \rightarrow 70.6^{\circ}\text{F}$  or colder.

**3.9** (a)  $N(25, 2.78)$ . (b)  $Z = 1.08 \rightarrow 0.1401$ .  
 (c) The answers are very close because only the units were changed. (The only reason why they differ at all is because  $28^\circ\text{C}$  is  $82.4^\circ\text{F}$ , not precisely  $83^\circ\text{F}$ .) (d) Since  $IQR = Q3 - Q1$ , we first need to find  $Q3$  and  $Q1$  and take the difference between the two. Remember that  $Q3$  is the  $75^{\text{th}}$  percentile and  $Q1$  is the  $25^{\text{th}}$  percentile of a distribution.  $Q1 = 23.13$ ,  $Q3 = 26.86$ ,  $IQR = 26.86 - 23.13 = 3.73$ .

**3.11** (a)  $Z = 0.67$ . (b)  $\mu = \$1650$ ,  $x = \$1800$ .  
 (c)  $0.67 = \frac{1800 - 1650}{\sigma} \rightarrow \sigma = \$223.88$ .

**3.13**  $Z = 1.56 \rightarrow 0.0594$ , i.e. 6%.

**3.15** (a)  $Z = 0.73 \rightarrow 0.2327$ . (b) If you are bidding on only one auction and set a low maximum bid price, someone will probably outbid you. If you set a high maximum bid price, you may win the auction but pay more than is necessary. If you are bidding on more than one auction, and you set your maximum bid price very low, you probably won't win any of the auctions. However, if the maximum bid price is even modestly high, you are likely to win multiple auctions. (c) An answer roughly equal to the 10th percentile would be reasonable. Regrettably, no percentile cutoff point guarantees beyond any possible event that you win at least one auction. However, you may pick a higher percentile if you want to be more sure of winning an auction. (d) Answers will vary a little but should correspond to the answer in part (c). We use the  $10^{\text{th}}$  percentile:  $Z = -1.28 \rightarrow \$69.80$ .

**3.17** (a) 70% of the data are within 1 standard deviation of the mean, 95% are within 2 and 100% are within 3 standard deviations of the mean. Therefore, we can say that the data approximately follow the 68-95-99.7% Rule. (b) The distribution is unimodal and symmetric. The superimposed normal curve seems to approximate the distribution pretty well. The points on the normal probability plot also seem to follow a straight line. There is one possible outlier on the lower end that is apparent in both graphs, but it is not too extreme. We can say that the distribution is nearly normal.

**??** (a) No. The cards are not independent. For example, if the first card is an ace of clubs, that implies the second card cannot be an ace

of clubs. Additionally, there are many possible categories, which would need to be simplified. (b) No. There are six events under consideration. The Bernoulli distribution allows for only two events or categories. Note that rolling a die could be a Bernoulli trial if we simply limit it to two events, e.g. rolling a 6 and not rolling a 6, though specifying such details would be necessary.

**??** (a)  $(1 - 0.471)^2 \times 0.471 = 0.1318$ .  
 (b)  $0.471^3 = 0.1045$ . (c)  $\mu = 1/0.471 = 2.12$ ,  $\sigma = \sqrt{2.38} = 1.54$ . (d)  $\mu = 1/0.30 = 3.33$ ,  $\sigma = 2.79$ . (e) When  $p$  is smaller, the event is rarer, meaning the expected number of trials before a success and the standard deviation of the waiting time are higher.

**??** (a)  $0.875^2 \times 0.125 = 0.096$ .  
 (b)  $\mu = 8$ ,  $\sigma = 7.48$ .

**??** (a) Binomial conditions are met: (1) Independent trials: In a random sample, whether or not one 18-20 year old has consumed alcohol does not depend on whether or not another one has. (2) Fixed number of trials:  $n = 10$ . (3) Only two outcomes at each trial: Consumed or did not consume alcohol. (4) Probability of a success is the same for each trial:  $p = 0.697$ . (b) 0.203. (c) 0.203. (d) 0.167. (e) 0.997.

**??** (a)  $\mu = 34.85$ ,  $\sigma = 3.25$  (b)  $Z = \frac{45 - 34.85}{3.25} = 3.12$ . 45 is more than 3 standard deviations away from the mean, we can assume that it is an unusual observation. Therefore yes, we would be surprised. (c) Using the normal approximation, 0.0009. With 0.5 correction, 0.0015.

**??** Want to find the probability that there will be 1,786 or more enrollees. Using the normal approximation: 0.0582. With a 0.5 correction: 0.0559.

**??** (a)  $1 - 0.75^3 = 0.5781$ . (b) 0.1406.  
 (c) 0.4219. (d)  $1 - 0.25^3 = 0.9844$ .

**??** (a) Geometric distribution: 0.109. (b) Binomial: 0.219. (c) Binomial: 0.137. (d)  $1 - 0.875^6 = 0.551$ . (e) Geometric: 0.084. (f) Using a binomial distribution with  $n = 6$  and  $p = 0.75$ , we see that  $\mu = 4.5$ ,  $\sigma = 1.06$ , and  $Z = 2.36$ . Since this is not within 2 SD, it may be considered unusual.

**??** 0 wins ( $-\$3$ ): 0.1458. 1 win ( $-\$1$ ): 0.3936. 2 wins ( $+\$1$ ): 0.3543. 3 wins ( $+\$3$ ): 0.1063.

?? (a)  $\frac{Anna}{1/5} \times \frac{Ben}{1/4} \times \frac{Carl}{1/3} \times \frac{Damian}{1/2} \times \frac{Eddy}{1/1} = \frac{1}{5!} = 1/120$ . (b) Since the probabilities must add to 1, there must be  $5! = 120$  possible orderings. (c)  $8! = 40,320$ .

?? (a) 0.0804. (b) 0.0322. (c) 0.0193.

?? (a) Negative binomial with  $n = 4$  and  $p = 0.55$ , where a success is defined here as a female student. The negative binomial setting is appropriate since the last trial is fixed but the order of the first 3 trials is unknown. (b) 0.1838. (c)  $\binom{3}{1} = 3$ . (d) In the binomial model there are

no restrictions on the outcome of the last trial. In the negative binomial model the last trial is fixed. Therefore we are interested in the number of ways of orderings of the other  $k - 1$  successes in the first  $n - 1$  trials.

?? (a) Poisson with  $\lambda = 75$ . (b)  $\mu = \lambda = 75$ ,  $\sigma = \sqrt{\lambda} = 8.66$ . (c)  $Z = -1.73$ . Since 60 is within 2 standard deviations of the mean, it would not generally be considered unusual. Note that we often use this rule of thumb even when the normal model does not apply. (d) Using Poisson with  $\lambda = 75$ : 0.0402.

## 4 Foundations for inference

**4.1** (a) Mean. Each student reports a numerical value: a number of hours. (b) Mean. Each student reports a number, which is a percentage, and we can average over these percentages. (c) Proportion. Each student reports Yes or No, so this is a categorical variable and we use a proportion. (d) Mean. Each student reports a number, which is a percentage like in part (b). (e) Proportion. Each student reports whether or not s/he expects to get a job, so this is a categorical variable and we use a proportion.

**4.3** (a) Mean: 13.65. Median: 14. (b) SD: 1.91. IQR:  $15 - 13 = 2$ . (c)  $Z_{16} = 1.23$ , which is not unusual since it is within 2 SD of the mean.  $Z_{18} = 2.28$ , which is generally considered unusual. (d) No. Point estimates that are based on samples only approximate the population parameter, and they vary from one sample to another. (e) We use the SE, which is  $1.91/\sqrt{100} = 0.191$  for this sample's mean.

**4.5** (a) We are building a distribution of sample statistics, in this case the sample mean. Such a distribution is called a sampling distribution. (b) Because we are dealing with the distribution of sample means, we need to check to see if the Central Limit Theorem applies. Our sample size is greater than 30, and we are told that random sampling is employed. With these conditions met, we expect that the distribution of the sample mean will be nearly normal and therefore symmetric. (c) Because we are dealing with a sampling distribution, we measure its variability with the standard error.  $SE = 18.2/\sqrt{45} = 2.713$ . (d) The sample means will be more variable with the smaller sample size.

**4.7** Recall that the general formula is

$$\text{point estimate} \pm Z^* \times SE$$

First, identify the three different values. The point estimate is 45%,  $Z^* = 1.96$  for a 95% confidence level, and  $SE = 1.2\%$ . Then, plug the values into the formula:

$$45\% \pm 1.96 \times 1.2\% \rightarrow (42.6\%, 47.4\%)$$

We are 95% confident that the proportion of US adults who live with one or more chronic conditions is between 42.6% and 47.4%.

**4.9** (a) False. Confidence intervals provide a range of plausible values, and sometimes the truth is missed. A 95% confidence interval "misses" about 5% of the time. (b) True. Notice that the description focuses on the true population value. (c) True. If we examine the 95% confidence interval computed in Exercise 4.9, we can see that 50% is not included in this interval. This means that in a hypothesis test, we would reject the null hypothesis that the proportion is 0.5. (d) False. The standard error describes the uncertainty in the overall estimate from natural fluctuations due to randomness, not the uncertainty corresponding to individuals' responses.

**4.11** (a) We are 95% confident that Americans spend an average of 1.38 to 1.92 hours per day relaxing or pursuing activities they enjoy. (b) Their confidence level must be higher as the width of the confidence interval increases as the confidence level increases. (c) The new margin of error will be smaller since as the sample size increases the standard error decreases, which will decrease the margin of error.

**4.13** (a) False. Provided the data distribution is not very strongly skewed ( $n = 64$  in this sample, so we can be slightly lenient with the skew), the sample mean will be nearly normal, allowing for the method normal approximation described. (b) False. Inference is made on the population parameter, not the point estimate. The point estimate is always in the confidence interval. (c) True. (d) False. The confidence interval is not about a sample mean. (e) False. To be more confident that we capture the parameter, we need a wider interval. Think about needing a bigger net to be more sure of catching a fish in a murky lake. (f) True. Optional explanation: This is true since the normal model was used to model the sample mean. The margin of error is half the width of the interval, and the sample mean is the midpoint of the interval. (g) False. In the calculation of the standard error, we divide the standard deviation by the square root of the sample size. To cut the SE (or margin of error) in half, we would need to sample  $2^2 = 4$  times the number of people in the initial sample.

**4.15** Independence: sample from  $< 10\%$  of population, and it is a random sample. We can assume that the students in this sample are independent of each other with respect to number of exclusive relationships they have been in. Notice that there are no students who have had no exclusive relationships in the sample, which suggests some student responses are likely missing (perhaps only positive values were reported). The sample size is at least 30. The skew is strong, but the sample is very large so this is not a concern. 90% CI: (2.97, 3.43). We are 90% confident that undergraduate students have been in 2.97 to 3.43 exclusive relationships, on average.

**4.17** (a)  $H_0 : \mu = 8$  (On average, New Yorkers sleep 8 hours a night.)

$H_A : \mu < 8$  (On average, New Yorkers sleep less than 8 hours a night.)

(b)  $H_0 : \mu = 15$  (The average amount of company time each employee spends not working is 15 minutes for March Madness.)

$H_A : \mu > 15$  (The average amount of company time each employee spends not working is greater than 15 minutes for March Madness.)

**4.19** The hypotheses should be about the pop-

ulation mean ( $\mu$ ), not the sample mean. The null hypothesis should have an equal sign and the alternative hypothesis should be about the null hypothesized value, not the observed sample mean. Correction:

$$H_0 : \mu = 10 \text{ hours}$$

$$H_A : \mu > 10 \text{ hours}$$

The one-sided test indicates that we are only interested in showing that 10 is an underestimate. Here the interest is in only one direction, so a one-sided test seems most appropriate. If we would also be interested if the data showed strong evidence that 10 was an overestimate, then the test should be two-sided.

**4.21** (a) This claim does not support since 3 hours (180 minutes) is not in the interval. (b) 2.2 hours (132 minutes) is in the 95% confidence interval, so we do not have evidence to say she is wrong. However, it would be more appropriate to use the point estimate of the sample. (c) A 99% confidence interval will be wider than a 95% confidence interval, meaning it would enclose this smaller interval. This means 132 minutes would be in the wider interval, and we would not reject her claim based on a 99% confidence level.

**4.23**  $H_0 : \mu = 130$ .  $H_A : \mu \neq 130$ .  $Z = 1.39 \rightarrow$  p-value = 0.1646, which is larger than  $\alpha = 0.05$ . The data do not provide convincing evidence that the true average calorie content in bags of potato chips is different than 130 calories.

**4.25** (a) Independence: The sample is random and 64 patients would almost certainly make up less than 10% of the ER residents. The sample size is at least 30. No information is provided about the skew. In practice, we would ask to see the data to check this condition, but here we will make the assumption that the skew is not very strong. (b)  $H_0 : \mu = 127$ .  $H_A : \mu \neq 127$ .  $Z = 2.15 \rightarrow$  p-value = 0.0316. Since the p-value is less than  $\alpha = 0.05$ , we reject  $H_0$ . The data provide convincing evidence that the average ER wait time has increased over the last year. (c) Yes, it would change. The p-value is greater than 0.01, meaning we would fail to reject  $H_0$  at  $\alpha = 0.01$ .

$$4.27 \quad Z = 1.65 = \frac{\bar{x} - 30}{10/\sqrt{70}} \rightarrow \bar{x} = 31.97.$$

**4.29** (a)  $H_0$ : Anti-depressants do not help symptoms of Fibromyalgia.  $H_A$ : Anti-depressants do treat symptoms of Fibromyalgia. Remark: Diana might also have taken special note if her symptoms got much worse, so a more scientific approach would have been to use a two-sided test. If you proposed a two-sided approach, your answers in (b) and (c) will be different. (b) Concluding that anti-depressants work for the treatment of Fibromyalgia symptoms when they actually do not. (c) Concluding that anti-depressants do not work for the treatment of Fibromyalgia symptoms when they actually do.

**4.31** (a) Scenario I is higher. Recall that a sample mean based on less data tends to be less accurate and have larger standard errors. (b) Scenario I is higher. The higher the confidence level, the higher the corresponding margin of error. (c) They are equal. The sample size does not affect the calculation of the p-value for a given Z-score. (d) Scenario I is higher. If the null hypothesis is harder to reject (lower  $\alpha$ ), then we are more likely to make a Type 2 Error when the alternative hypothesis is true.

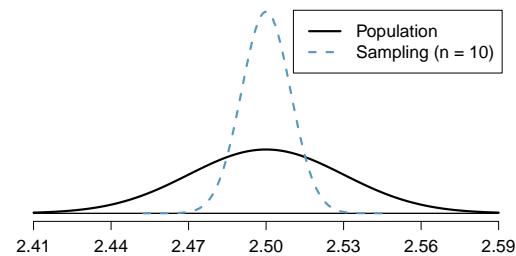
**4.33** (a) The distribution is unimodal and strongly right skewed with a median between 5 and 10 years old. Ages range from 0 to slightly over 50 years old, and the middle 50% of the distribution is roughly between 5 and 15 years old. There are potential outliers on the higher end. (b) When the sample size is small, the sampling distribution is right skewed, just like the population distribution. As the sample size increases, the sampling distribution gets more unimodal, symmetric, and approaches normality. The variability also decreases. This is consistent with the Central Limit Theorem. (c)  $n = 5$ :  $\mu_{\bar{x}} = 10.44$ ,  $\sigma_{\bar{x}} = 4.11$ ;  $n = 30$ :  $\mu_{\bar{x}} = 10.44$ ,  $\sigma_{\bar{x}} = 1.68$ ;  $n = 100$ :  $\mu_{\bar{x}} = 10.44$ ,  $\sigma_{\bar{x}} = 0.92$ . The centers of the sampling distributions shown in part (b) appear to be around 10. It is difficult to estimate the standard deviation for the sampling distribution when  $n = 5$  from the histogram (since the distribution is somewhat skewed). If 1.68 is a plausible estimate for the standard deviation of the sampling distribution when  $n = 30$ , then using the 68-95-99.7% Rule, we would expect the values to range roughly between  $10.44 \pm 3 * 1.68 = (5.4, 15.48)$ , which seems to be the case. Similarly, when  $n = 100$ , we would expect the values to range roughly be-

tween  $10.44 \pm 3 * 0.92 = (7.68, 13.2)$ , which also seems to be the case.

**4.35** (a) Right skewed. There is a long tail on the higher end of the distribution but a much shorter tail on the lower end. (b) Less than, as the median would be less than the mean in a right skewed distribution. (c) We should not. (d) Even though the population distribution is not normal, the conditions for inference are reasonably satisfied, with the possible exception of skew. If the skew isn't very strong (we should ask to see the data), then we can use the Central Limit Theorem to estimate this probability. For now, we'll assume the skew isn't very strong, though the description suggests it is at least moderate to strong. Use  $N(1.3, SD_{\bar{x}} = 0.3/\sqrt{60})$ :  $Z = 2.58 \rightarrow 0.0049$ . (e) It would decrease it by a factor of  $1/\sqrt{2}$ .

**4.37** The centers are the same in each plot, and each data set is from a nearly normal distribution, though the histograms may not look very normal since each represents only 100 data points. The only way to tell which plot corresponds to which scenario is to examine the variability of each distribution. Plot B is the most variable, followed by Plot A, then Plot C. This means Plot B will correspond to the original data, Plot A to the sample means with size 5, and Plot C to the sample means with size 25.

**4.39** (a)  $Z = -3.33 \rightarrow 0.0004$ . (b) The population SD is known and the data are nearly normal, so the sample mean will be nearly normal with distribution  $N(\mu, \sigma/\sqrt{n})$ , i.e.  $N(2.5, 0.0095)$ . (c)  $Z = -10.54 \rightarrow \approx 0$ . (d) See below:



(e) We could not estimate (a) without a nearly normal population distribution. We also could not estimate (c) since the sample size is not sufficient to yield a nearly normal sampling distribution if the population distribution is not nearly normal.

**4.41** (a) We cannot use the normal model for this calculation, but we can use the histogram. About 500 songs are shown to be longer than 5 minutes, so the probability is about  $500/3000 = 0.167$ . (b) Two different answers are reasonable. *Option*<sup>1</sup> Since the population distribution is only slightly skewed to the right, even a small sample size will yield a nearly normal sampling distribution. We also know that the songs are sampled randomly and the sample size is less than 10% of the population, so the length of one song in the sample is independent of another. We are looking for the probability that the total length of 15 songs is more than 60 minutes, which means that the average song should last at least  $60/15 = 4$  minutes. Using  $SD_{\bar{x}} = 1.63/\sqrt{15}$ ,  $Z = 1.31 \rightarrow 0.0951$ . *Option*<sup>2</sup> Since the population distribution is not normal, a small sample size may not be sufficient to yield a nearly normal sampling distribution. Therefore, we cannot estimate the probability using the tools we have learned so far. (c) We can now be confident that the conditions are satisfied.  $Z = 0.92 \rightarrow 0.1788$ .

**??** (a)  $H_0 : \mu_{2009} = \mu_{2004}$ .  $H_A : \mu_{2009} \neq \mu_{2004}$ . (b)  $\bar{x}_{2009} - \bar{x}_{2004} = -3.6$  spam emails per day. (c) The null hypothesis was not rejected, and the data do not provide convincing evidence that the true average number of spam emails per day in years 2004 and 2009 are different. The observed difference is about what we might expect from sampling variability alone. (d) Yes, since the hypothesis of no difference was not rejected in part (c).

**??** (a)  $H_0 : p_{2009} = p_{2004}$ .  $H_A : p_{2009} \neq p_{2004}$ . (b) -7%. (c) The null hypothesis was rejected. The data provide strong evidence that the true proportion of those who once a month or less frequently delete their spam email was higher in 2004 than in 2009. The difference is so large that it cannot easily be explained as being due to chance. (d) No, since the null difference, 0, was rejected in part (c).

**??** True. If the sample size is large, then the standard error will be small, meaning even relatively small differences between the null value and point estimate can be statistically significant.

## 5 Inference for numerical data

**5.1** (a)  $df = 6 - 1 = 5$ ,  $t_5^* = 2.02$  (column with two tails of 0.10, row with  $df = 5$ ). (b)  $df = 21 - 1 = 20$ ,  $t_{20}^* = 2.53$  (column with two tails of 0.02, row with  $df = 20$ ). (c)  $df = 28$ ,  $t_{28}^* = 2.05$ . (d)  $df = 11$ ,  $t_{11}^* = 3.11$ .

**5.3** (a) between 0.025 and 0.05 (b) less than 0.005 (c) greater than 0.2 (d) between 0.01 and 0.025

**5.5** The mean is the midpoint:  $\bar{x} = 20$ . Identify the margin of error:  $ME = 1.015$ , then use  $t_{35}^* = 2.03$  and  $SE = s/\sqrt{n}$  in the formula for margin of error to identify  $s = 3$ .

**5.7** (a)  $H_0: \mu = 8$  (New Yorkers sleep 8 hrs per night on average.)  $H_A: \mu < 8$  (New Yorkers sleep less than 8 hrs per night on average.) (b) Independence: The sample is random and

from less than 10% of New Yorkers. The sample is small, so we will use a *t*-distribution. For this size sample, slight skew is acceptable, and the min/max suggest there is not much skew in the data.  $T = -1.75$ .  $df = 25 - 1 = 24$ . (c)  $0.025 < p\text{-value} < 0.05$ . If in fact the true population mean of the amount New Yorkers sleep per night was 8 hours, the probability of getting a random sample of 25 New Yorkers where the average amount of sleep is 7.73 hrs per night or less is between 0.025 and 0.05. (d) Since  $p\text{-value} < 0.05$ , reject  $H_0$ . The data provide strong evidence that New Yorkers sleep less than 8 hours per night on average. (e) No, as we rejected  $H_0$ .

**5.9**  $t_{19}^*$  is 1.73 for a one-tail. We want the lower tail, so set -1.73 equal to the T-score, then solve for  $\bar{x}$ : 56.91.

**5.11** (a) We will conduct a 1-sample  $t$ -test.  $H_0: \mu = 5$ .  $H_A: \mu < 5$ . We'll use  $\alpha = 0.05$ . This is a random sample, so the observations are independent. To proceed, we assume the distribution of years of piano lessons is approximately normal.  $SE = 2.2/\sqrt{20} = 0.4919$ . The test statistic is  $T = (4.6 - 5)/SE = -0.81$ .  $df = 20 - 1 = 19$ . The one-tail p-value is about 0.21, which is bigger than  $\alpha = 0.05$ , so we do not reject  $H_0$ . That is, we do not have sufficiently strong evidence to reject Georgianna's claim.

(b) Using  $SE = 0.4919$  and  $t_{df=19}^* = 2.093$ , the confidence interval is  $(3.57, 5.63)$ . We are 95% confident that the average number of years a child takes piano lessons in this city is 3.57 to 5.63 years.

(c) They agree, since we did not reject the null hypothesis and the null value of 5 was in the  $t$ -interval.

**5.13** If the sample is large, then the margin of error will be about  $1.96 \times 100/\sqrt{n}$ . We want this value to be less than 10, which leads to  $n \geq 384.16$ , meaning we need a sample size of at least 385 (round up for sample size calculations!).

**5.15** (a) Two-sided, we are evaluating a difference, not in a particular direction. (b) Paired, data are recorded in the same cities at two different time points. The temperature in a city at one point is not independent of the temperature in the same city at another time point. (c)  $t$ -test, sample is small and population standard deviation is unknown.

**5.17** (a) Since it's the same students at the beginning and the end of the semester, there is a pairing between the datasets, for a given student their beginning and end of semester grades are dependent. (b) Since the subjects were sampled randomly, each observation in the men's group does not have a special correspondence with exactly one observation in the other (women's) group. (c) Since it's the same subjects at the beginning and the end of the study, there is a pairing between the datasets, for a subject stu-

dent their beginning and end of semester artery thickness are dependent. (d) Since it's the same subjects at the beginning and the end of the study, there is a pairing between the datasets, for a subject student their beginning and end of semester weights are dependent.

**5.19** (a) For each observation in one data set, there is exactly one specially-corresponding observation in the other data set for the same geographic location. The data are paired. (b)  $H_0: \mu_{diff} = 0$  (There is no difference in average daily high temperature between January 1, 1968 and January 1, 2008 in the continental US.)  $H_A: \mu_{diff} > 0$  (Average daily high temperature in January 1, 1968 was lower than average daily high temperature in January, 2008 in the continental US.) If you chose a two-sided test, that would also be acceptable. If this is the case, note that your p-value will be a little bigger than what is reported here in part (d). (c) Locations are random and represent less than 10% of all possible locations in the US. The sample size is at least 30. We are not given the distribution to check the skew. In practice, we would ask to see the data to check this condition, but here we will move forward under the assumption that it is not strongly skewed. (d)  $T_{50} \approx 1.60 \rightarrow 0.05 < \text{p-value} < 0.10$ . (e) Since the p-value  $> \alpha$  (since not given use 0.05), fail to reject  $H_0$ . The data do not provide strong evidence of temperature warming in the continental US. However it should be noted that the p-value is very close to 0.05. (f) Type 2 Error, since we may have incorrectly failed to reject  $H_0$ . There may be an increase, but we were unable to detect it. (g) Yes, since we failed to reject  $H_0$ , which had a null value of 0.

**5.21** (a)  $(-0.05, 2.25)$ . (b) We are 90% confident that the average daily high on January 1, 2008 in the continental US was 0.05 degrees lower to 2.25 degrees higher than the average daily high on January 1, 1968. (c) No, since 0 is included in the interval.

**5.23** (a) Each of the 36 mothers is related to exactly one of the 36 fathers (and vice-versa), so there is a special correspondence between the mothers and fathers. (b)  $H_0 : \mu_{diff} = 0$ .  $H_A : \mu_{diff} \neq 0$ . Independence: random sample from less than 10% of population. Sample size of at least 30. The skew of the differences is, at worst, slight.  $T_{35} = 2.72 \rightarrow p\text{-value} = 0.01$ . Since  $p\text{-value} < 0.05$ , reject  $H_0$ . The data provide strong evidence that the average IQ scores of mothers and fathers of gifted children are different, and the data indicate that mothers' scores are higher than fathers' scores for the parents of gifted children.

**5.25** No, he should not move forward with the test since the distributions of total personal income are very strongly skewed. When sample sizes are large, we can be a bit lenient with skew. However, such strong skew observed in this exercise would require somewhat large sample sizes, somewhat higher than 30.

**5.27** (a) These data are paired. For example, the Friday the 13th in say, September 1991, would probably be more similar to the Friday the 6th in September 1991 than to Friday the 6th in another month or year. (b) Let  $\mu_{diff} = \mu_{sixth} - \mu_{thirteenth}$ .  $H_0 : \mu_{diff} = 0$ .  $H_A : \mu_{diff} \neq 0$ . (c) Independence: The months selected are not random. However, if we think these dates are roughly equivalent to a simple random sample of all such Friday 6th/13th date pairs, then independence is reasonable. To proceed, we must make this strong assumption, though we should note this assumption in any reported results. Normality: With fewer than 10 observations, we would need to use the  $t$ -distribution to model the sample mean. The normal probability plot of the differences shows an approximately straight line. There isn't a clear reason why this distribution would be skewed, and since the normal quantile plot looks reasonable, we can mark this condition as reasonably satisfied. (d)  $T = 4.94$  for  $df = 10 - 1 = 9 \rightarrow p\text{-value} < 0.01$ . (e) Since  $p\text{-value} < 0.05$ , reject  $H_0$ . The data provide strong evidence that the average number of cars at the intersection is higher on Friday the 6<sup>th</sup> than on Friday the 13<sup>th</sup>. (We might believe this intersection is representative of all roads, i.e. there is higher traffic on Friday the 6<sup>th</sup> relative to Friday the 13<sup>th</sup>.

However, we should be cautious of the required assumption for such a generalization.) (f) If the average number of cars passing the intersection actually was the same on Friday the 6<sup>th</sup> and 13<sup>th</sup>, then the probability that we would observe a test statistic so far from zero is less than 0.01. (g) We might have made a Type 1 Error, i.e. incorrectly rejected the null hypothesis.

**5.29** (a)  $H_0 : \mu_{diff} = 0$ .  $H_A : \mu_{diff} \neq 0$ .  $T = -2.71$ .  $df = 5$ .  $0.02 < p\text{-value} < 0.05$ . Since  $p\text{-value} < 0.05$ , reject  $H_0$ . The data provide strong evidence that the average number of traffic accident related emergency room admissions are different between Friday the 6<sup>th</sup> and Friday the 13<sup>th</sup>. Furthermore, the data indicate that the direction of that difference is that accidents are lower on Friday the 6<sup>th</sup> relative to Friday the 13<sup>th</sup>. (b) (-6.49, -0.17). (c) This is an observational study, not an experiment, so we cannot so easily infer a causal intervention implied by this statement. It is true that there is a difference. However, for example, this does not mean that a responsible adult going out on Friday the 13<sup>th</sup> has a higher chance of harm than on any other night.

**5.31** (a) Chicken fed linseed weighed an average of 218.75 grams while those fed horsebean weighed an average of 160.20 grams. Both distributions are relatively symmetric with no apparent outliers. There is more variability in the weights of chicken fed linseed. (b)  $H_0 : \mu_{ls} = \mu_{hb}$ .  $H_A : \mu_{ls} \neq \mu_{hb}$ . We leave the conditions to you to consider.  $T = 3.02$ ,  $df = \min(11, 9) = 9 \rightarrow 0.01 < p\text{-value} < 0.02$ . Since  $p\text{-value} < 0.05$ , reject  $H_0$ . The data provide strong evidence that there is a significant difference between the average weights of chickens that were fed linseed and horsebean. (c) Type 1 Error, since we rejected  $H_0$ . (d) Yes, since  $p\text{-value} > 0.01$ , we would have failed to reject  $H_0$ .

**5.33**  $H_0 : \mu_C = \mu_S$ .  $H_A : \mu_C \neq \mu_S$ .  $T = 3.27$ ,  $df = 11 \rightarrow p\text{-value} < 0.01$ . Since  $p\text{-value} < 0.05$ , reject  $H_0$ . The data provide strong evidence that the average weight of chickens that were fed casein is different than the average weight of chickens that were fed soybean (with weights from casein being higher). Since this is a randomized experiment, the observed difference can be attributed to the diet.

**5.35**  $H_0 : \mu_T = \mu_C$ .  $H_A : \mu_T \neq \mu_C$ .  $T = 2.24$ ,  $df = 21 \rightarrow 0.02 < p\text{-value} < 0.05$ . Since  $p\text{-value} < 0.05$ , reject  $H_0$ . The data provide strong evidence that the average food consumption by the patients in the treatment and control groups are different. Furthermore, the data indicate patients in the distracted eating (treatment) group consume more food than patients in the control group.

**5.37** Let  $\mu_{diff} = \mu_{pre} - \mu_{post}$ .  $H_0 : \mu_{diff} = 0$ : Treatment has no effect.  $H_A : \mu_{diff} > 0$ : Treatment is effective in reducing P.D.T. scores, the average pre-treatment score is higher than the average post-treatment score. Note that the reported values are pre minus post, so we are looking for a positive difference, which would correspond to a reduction in the P.D.T. score. Conditions are checked as follows. Independence: The subjects are randomly assigned to treatments, so the patients in each group are independent. All three sample sizes are smaller than 30, so we use  $t$ -tests. Distributions of differences are somewhat skewed. The sample sizes are small, so we cannot reliably relax this assumption. (We will proceed, but we would not report the results of this specific analysis, at least for treatment group 1.) For all three groups:  $df = 13$ .  $T_1 = 1.89$  ( $0.025 < p\text{-value} < 0.05$ ),  $T_2 = 1.35$  ( $p\text{-value} = 0.10$ ),  $T_3 = -1.40$  ( $p\text{-value} > 0.10$ ). The only significant test reduction is found in Treatment 1, however, we had earlier noted that this result might not be reliable due to the skew in the distribution. Note that the calculation of the  $p$ -value for Treatment 3 was unnecessary: the sample mean indicated a increase in P.D.T. scores under this treatment (as opposed to a decrease, which was the result of interest). That is, we could tell without formally completing the hypothesis test that the  $p$ -value would be large for this treatment group.

**5.39** Difference we care about: 40. Single tail of 90%:  $1.28 \times SE$ . Rejection region bounds:  $\pm 1.96 \times SE$  (if 5% significance level). Setting  $3.24 \times SE = 40$ , subbing in  $SE = \sqrt{\frac{94^2}{n}} + \frac{94^2}{n}$ , and solving for the sample size  $n$  gives 116 plots of land for each fertilizer.

**5.41** Alternative.

**5.43**  $H_0: \mu_1 = \mu_2 = \dots = \mu_6$ .  $H_A$ : The average weight varies across some (or all) groups. Independence: Chicks are randomly assigned to feed types (presumably kept separate from one another), therefore independence of observations is reasonable. Approx. normal: the distributions of weights within each feed type appear to be fairly symmetric. Constant variance: Based on the side-by-side box plots, the constant variance assumption appears to be reasonable. There are differences in the actual computed standard deviations, but these might be due to chance as these are quite small samples.  $F_{5,65} = 15.36$  and the  $p$ -value is approximately 0. With such a small  $p$ -value, we reject  $H_0$ . The data provide convincing evidence that the average weight of chicks varies across some (or all) feed supplement groups.

**5.45** (a)  $H_0$ : The population mean of MET for each group is equal to the others.  $H_A$ : At least one pair of means is different. (b) Independence: We don't have any information on how the data were collected, so we cannot assess independence. To proceed, we must assume the subjects in each group are independent. In practice, we would inquire for more details. Approx. normal: The data are bound below by zero and the standard deviations are larger than the means, indicating very strong skew. However, since the sample sizes are extremely large, even extreme skew is acceptable. Constant variance: This condition is sufficiently met, as the standard deviations are reasonably consistent across groups. (c) See below, with the last column omitted:

	Df	Sum Sq	Mean Sq	F value
coffee	4	10508	2627	5.2
Residuals	50734	25564819	504	
Total	50738	25575327		

(d) Since  $p$ -value is very small, reject  $H_0$ . The data provide convincing evidence that the average MET differs between at least one pair of groups.

**5.47** (a)  $H_0$ : Average GPA is the same for all majors.  $H_A$ : At least one pair of means are different. (b) Since p-value > 0.05, fail to reject  $H_0$ . The data do not provide convincing evidence of a difference between the average GPAs across three groups of majors. (c) The total degrees of freedom is  $195 + 2 = 197$ , so the sample size is  $197 + 1 = 198$ .

**5.49** (a) False. As the number of groups increases, so does the number of comparisons and hence the modified significance level decreases. (b) True. (c) True. (d) False. We need observations to be independent regardless of sample size.

**??** (a)  $H_0$ : Average score difference is the same for all treatments.  $H_A$ : At least one pair of means are different. (b) We should check conditions. If we look back to the earlier exercise, we will see that the patients were randomized, so independence is satisfied. There are some minor concerns about skew, especially with the third group, though this may be acceptable. The

standard deviations across the groups are reasonably similar. Since the p-value is less than 0.05, reject  $H_0$ . The data provide convincing evidence of a difference between the average reduction in score among treatments. (c) We determined that at least two means are different in part (b), so we now conduct  $K = 3 \times 2/2 = 3$  pairwise  $t$ -tests that each use  $\alpha = 0.05/3 = 0.0167$  for a significance level. Use the following hypotheses for each pairwise test.  $H_0$ : The two means are equal.  $H_A$ : The two means are different. The sample sizes are equal and we use the pooled SD, so we can compute  $SE = 3.7$  with the pooled  $df = 39$ . The p-value only for Trmt 1 vs. Trmt 3 may be statistically significant:  $0.01 < \text{p-value} < 0.02$ . Since we cannot tell, we should use a computer to get the p-value, 0.015, which is statistically significant for the adjusted significance level. That is, we have identified Treatment 1 and Treatment 3 as having different effects. Checking the other two comparisons, the differences are not statistically significant.

## 6 Inference for categorical data

**??** (a) False. Doesn't satisfy success-failure condition. (b) True. The success-failure condition is not satisfied. In most samples we would expect  $\hat{p}$  to be close to 0.08, the true population proportion. While  $\hat{p}$  can be much above 0.08, it is bound below by 0, suggesting it would take on a right skewed shape. Plotting the sampling distribution would confirm this suspicion. (c) False.  $SE_{\hat{p}} = 0.0243$ , and  $\hat{p} = 0.12$  is only  $\frac{0.12 - 0.08}{0.0243} = 1.65$  SEs away from the mean, which would not be considered unusual. (d) True.  $\hat{p} = 0.12$  is 2.32 standard errors away from the mean, which is often considered unusual. (e) False. Decreases the SE by a factor of  $1/\sqrt{2}$ .

**??** (a) True. See the reasoning of 6.1(b). (b) True. We take the square root of the sample

size in the SE formula. (c) True. The independence and success-failure conditions are satisfied. (d) True. The independence and success-failure conditions are satisfied.

**??** (a) False. A confidence interval is constructed to estimate the population proportion, not the sample proportion. (b) True. 95% CI:  $70\% \pm 8\%$ . (c) True. By the definition of the confidence level. (d) True. Quadrupling the sample size decreases the SE and ME by a factor of  $1/\sqrt{4}$ . (e) True. The 95% CI is entirely above 50%.

**??** With a random sample from < 10% of the population, independence is satisfied. The success-failure condition is also satisfied.  $ME = z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = 1.96 \sqrt{\frac{0.56 \times 0.44}{600}} = 0.0397 \approx 4\%$

?? (a) Proportion of graduates from this university who found a job within one year of graduating.  $\hat{p} = 348/400 = 0.87$ . (b) This is a random sample from less than 10% of the population, so the observations are independent. Success-failure condition is satisfied: 348 successes, 52 failures, both well above 10. (c) (0.8371, 0.9029). We are 95% confident that approximately 84% to 90% of graduates from this university found a job within one year of completing their undergraduate degree. (d) 95% of such random samples would produce a 95% confidence interval that includes the true proportion of students at this university who found a job within one year of graduating from college. (e) (0.8267, 0.9133). Similar interpretation as before. (f) 99% CI is wider, as we are more confident that the true proportion is within the interval and so need to cover a wider range.

?? (a) No. The sample only represents students who took the SAT, and this was also an online survey. (b) (0.5289, 0.5711). We are 90% confident that 53% to 57% of high school seniors who took the SAT are fairly certain that they will participate in a study abroad program in college. (c) 90% of such random samples would produce a 90% confidence interval that includes the true proportion. (d) Yes. The interval lies entirely above 50%.

?? (a) This is an appropriate setting for a hypothesis test.  $H_0 : p = 0.50$ .  $H_A : p > 0.50$ . Both independence and the success-failure condition are satisfied.  $Z = 1.12 \rightarrow \text{p-value} = 0.1314$ . Since the p-value  $> \alpha = 0.05$ , we fail to reject  $H_0$ . The data do not provide strong evidence that more than half of all Independents oppose the public option plan. (b) Yes, since we did not reject  $H_0$  in part (a).

?? (a)  $H_0 : p = 0.38$ .  $H_A : p \neq 0.38$ . Independence (random sample, < 10% of population) and the success-failure condition are satisfied.  $Z = -20.5 \rightarrow \text{p-value} \approx 0$ . Since the p-value is very small, we reject  $H_0$ . The data provide strong evidence that the proportion of Americans who only use their cell phones to access

the internet is different than the Chinese proportion of 38%, and the data indicate that the proportion is lower in the US. (b) If in fact 38% of Americans used their cell phones as a primary access point to the internet, the probability of obtaining a random sample of 2,254 Americans where 17% or less or 59% or more use their only their cell phones to access the internet would be approximately 0. (c) (0.1545, 0.1855). We are 95% confident that approximately 15.5% to 18.6% of all Americans primarily use their cell phones to browse the internet.

?? (a)  $H_0 : p = 0.5$ .  $H_A : p > 0.5$ . Independence (random sample, < 10% of population) is satisfied, as is the success-failure conditions (using  $p_0 = 0.5$ , we expect 40 successes and 40 failures).  $Z = 2.91 \rightarrow \text{p-value} = 0.0018$ . Since the p-value  $< 0.05$ , we reject the null hypothesis. The data provide strong evidence that the rate of correctly identifying a soda for these people is significantly better than just by random guessing. (b) If in fact people cannot tell the difference between diet and regular soda and they randomly guess, the probability of getting a random sample of 80 people where 53 or more identify a soda correctly would be 0.0018.

?? (a) Independence is satisfied (random sample from < 10% of the population), as is the success-failure condition (40 smokers, 160 non-smokers). The 95% CI: (0.145, 0.255). We are 95% confident that 14.5% to 25.5% of all students at this university smoke. (b) We want  $z^*SE$  to be no larger than 0.02 for a 95% confidence level. We use  $z^* = 1.96$  and plug in the point estimate  $\hat{p} = 0.2$  within the SE formula:  $1.96\sqrt{0.2(1-0.2)/n} \leq 0.02$ . The sample size  $n$  should be at least 1,537.

?? The margin of error, which is computed as  $z^*SE$ , must be smaller than 0.01 for a 90% confidence level. We use  $z^* = 1.65$  for a 90% confidence level, and we can use the point estimate  $\hat{p} = 0.52$  in the formula for  $SE$ .  $1.65\sqrt{0.52(1-0.52)/n} \leq 0.01$ . Therefore, the sample size  $n$  must be at least 6,796.

**??** This is not a randomized experiment, and it is unclear whether people would be affected by the behavior of their peers. That is, independence may not hold. Additionally, there are only 5 interventions under the provocative scenario, so the success-failure condition does not hold. Even if we consider a hypothesis test where we pool the proportions, the success-failure condition will not be satisfied. Since one condition is questionable and the other is not satisfied, the difference in sample proportions will not follow a nearly normal distribution.

**??** (a) False. The entire confidence interval is above 0. (b) True. (c) True. (d) True. (e) False. It is simply the negated and reordered values: (-0.06,-0.02).

**??** (a) (0.23, 0.33). We are 95% confident that the proportion of Democrats who support the plan is 23% to 33% higher than the proportion of Independents who do. (b) True.

**??** (a) College grads: 23.7%. Non-college grads: 33.7%. (b) Let  $p_{CG}$  and  $p_{NCG}$  represent the proportion of college graduates and non-college graduates who responded “do not know”.  $H_0 : p_{CG} = p_{NCG}$ .  $H_A : p_{CG} \neq p_{NCG}$ . Independence is satisfied (random sample, < 10% of the population), and the success-failure condition, which we would check using the pooled proportion ( $\hat{p} = 235/827 = 0.284$ ), is also satisfied.  $Z = -3.18 \rightarrow$  p-value = 0.0014. Since the p-value is very small, we reject  $H_0$ . The data provide strong evidence that the proportion of college graduates who do not have an opinion on this issue is different than that of non-college graduates. The data also indicate that fewer college grads say they “do not know” than non-college grads (i.e. the data indicate the direction after we reject  $H_0$ ).

**??** (a) College grads: 35.2%. Non-college grads: 33.9%. (b) Let  $p_{CG}$  and  $p_{NCG}$  represent the proportion of college graduates and non-college grads who support offshore drilling.  $H_0 : p_{CG} = p_{NCG}$ .  $H_A : p_{CG} \neq p_{NCG}$ . Independence is satisfied (random sample, < 10% of the population), and the success-failure condition, which we would check using the pooled proportion ( $\hat{p} = 286/827 = 0.346$ ), is also satisfied.  $Z = 0.39 \rightarrow$  p-value = 0.6966. Since the p-value  $> \alpha$  (0.05), we fail to reject  $H_0$ . The

data do not provide strong evidence of a difference between the proportions of college graduates and non-college graduates who support offshore drilling in California.

**??** Subscript  $C$  means control group. Subscript  $T$  means truck drivers.  $H_0 : p_C = p_T$ .  $H_A : p_C \neq p_T$ . Independence is satisfied (random samples, < 10% of the population), as is the success-failure condition, which we would check using the pooled proportion ( $\hat{p} = 70/495 = 0.141$ ).  $Z = -1.65 \rightarrow$  p-value = 0.0989. Since the p-value is high (default to  $\alpha = 0.05$ ), we fail to reject  $H_0$ . The data do not provide strong evidence that the rates of sleep deprivation are different for non-transportation workers and truck drivers.

**??** (a) Summary of the study:

		<i>Virol. failure</i>		Total
		Yes	No	
<i>Treatment</i>	Nevaripine	26	94	120
	Lopinavir	10	110	120
Total		36	204	240

(b)  $H_0 : p_N = p_L$ . There is no difference in virologic failure rates between the Nevaripine and Lopinavir groups.  $H_A : p_N \neq p_L$ . There is some difference in virologic failure rates between the Nevaripine and Lopinavir groups. (c) Random assignment was used, so the observations in each group are independent. If the patients in the study are representative of those in the general population (something impossible to check with the given information), then we can also confidently generalize the findings to the population. The success-failure condition, which we would check using the pooled proportion ( $\hat{p} = 36/240 = 0.15$ ), is satisfied.  $Z = 2.89 \rightarrow$  p-value = 0.0039. Since the p-value is low, we reject  $H_0$ . There is strong evidence of a difference in virologic failure rates between the Nevaripine and Lopinavir groups. Treatment and virologic failure do not appear to be independent.

**??** No. The samples at the beginning and at the end of the semester are not independent since the survey is conducted on the same students.

**??** (a) False. The chi-square distribution has one parameter called degrees of freedom. (b) True. (c) True. (d) False. As the degrees of freedom increases, the shape of the chi-square distribution becomes more symmetric.

?? (a)  $H_0$ : The distribution of the format of the book used by the students follows the professor's predictions.  $H_A$ : The distribution of the format of the book used by the students does not follow the professor's predictions. (b)  $E_{\text{hard copy}} = 126 \times 0.60 = 75.6$ .  $E_{\text{print}} = 126 \times 0.25 = 31.5$ .  $E_{\text{online}} = 126 \times 0.15 = 18.9$ . (c) Independence: The sample is not random. However, if the professor has reason to believe that the proportions are stable from one term to the next and students are not affecting each other's study habits, independence is probably reasonable. Sample size: All expected counts are at least 5. (d)  $\chi^2 = 2.32$ ,  $df = 2$ , p-value > 0.3. (e) Since the p-value is large, we fail to reject  $H_0$ . The data do not provide strong evidence indicating the professor's predictions were statistically inaccurate.

?? Use a chi-squared goodness of fit test.  $H_0$ : Each option is equally likely.  $H_A$ : Some options are preferred over others. Total sample size: 99. Expected counts:  $(1/3) * 99 = 33$  for each option. These are all above 5, so conditions are satisfied.  $df = 3 - 1 = 2$  and  $\chi^2 = \frac{(43-33)^2}{33} + \frac{(21-33)^2}{33} + \frac{(35-33)^2}{33} = 7.52 \rightarrow 0.02 < \text{p-value} < 0.05$ . Since the p-value is less than 5%, we reject  $H_0$ . The data provide convincing evidence that some options are preferred over others.

?? (a) Two-way table:

Treatment	Quit		Total
	Yes	No	
Patch + support group	40	110	150
Only patch	30	120	150
Total	70	230	300

$$(b-i) E_{\text{row}_1,\text{col}_1} = \frac{(\text{row 1 total}) \times (\text{col 1 total})}{\text{table total}} = 35.$$

This is lower than the observed value.

$$(b-ii) E_{\text{row}_2,\text{col}_2} = \frac{(\text{row 2 total}) \times (\text{col 2 total})}{\text{table total}} = 115.$$

This is lower than the observed value.

??  $H_0$ : The opinion of college grads and non-grads is not different on the topic of drilling for oil and natural gas off the coast of California.  $H_A$ : Opinions regarding the drilling for oil and natural gas off the coast of California has an association with earning a college degree.

$$E_{\text{row 1,col 1}} = 151.5 \quad E_{\text{row 1,col 2}} = 134.5$$

$$E_{\text{row 2,col 1}} = 162.1 \quad E_{\text{row 2,col 2}} = 143.9$$

$$E_{\text{row 3,col 1}} = 124.5 \quad E_{\text{row 3,col 2}} = 110.5$$

Independence: The samples are both random, unrelated, and from less than 10% of the population, so independence between observations is reasonable. Sample size: All expected counts are at least 5.  $\chi^2 = 11.47$ ,  $df = 2 \rightarrow 0.001 < \text{p-value} < 0.005$ . Since the p-value <  $\alpha$ , we reject  $H_0$ . There is strong evidence that there is an association between support for off-shore drilling and having a college degree.

?? (a)  $H_0$ : The age of Los Angeles residents is independent of shipping carrier preference variable.  $H_A$ : The age of Los Angeles residents is associated with the shipping carrier preference variable. (b) The conditions are not satisfied since some expected counts are below 5.

?? No. For a confidence interval, we check the success-failure condition using the data, and there are only 9 respondents who said bullying is no problem at all.

?? (a)  $H_0 : p = 0.69$ .  $H_A : p \neq 0.69$ . (b)  $\hat{p} = \frac{17}{30} = 0.57$ . (c) The success-failure condition is not satisfied; note that it is appropriate to use the null value ( $p_0 = 0.69$ ) to compute the expected number of successes and failures. (d) Answers may vary. Each student can be represented with a card. Take 100 cards, 69 black cards representing those who follow the news about Egypt and 31 red cards representing those who do not. Shuffle the cards and draw with replacement (shuffling each time in between draws) 30 cards representing the 30 high school students. Calculate the proportion of black cards in this sample,  $\hat{p}_{\text{sim}}$ , i.e. the proportion of those who follow the news in the simulation. Repeat this many times (e.g. 10,000 times) and plot the resulting sample proportions. The p-value will be two times the proportion of simulations where  $\hat{p}_{\text{sim}} \leq 0.57$ . (Note: we would generally use a computer to perform these simulations.) (e) The p-value is about  $0.001 + 0.005 + 0.020 + 0.035 + 0.075 = 0.136$ , meaning the two-sided p-value is about 0.272. Your p-value may vary slightly since it is based on a visual estimate. Since the p-value is greater than 0.05, we fail to reject  $H_0$ . The data do not provide strong evidence that the proportion of high school students who followed the news about Egypt is different than the proportion of American adults who did.

**??** The subscript  $pr$  corresponds to provocative and  $con$  to conservative. (a)  $H_0 : p_{pr} = p_{con}$ .  $H_A : p_{pr} \neq p_{con}$ . (b) -0.35. (c) The left tail for the p-value is calculated by adding up the two left bins:  $0.005 + 0.015 = 0.02$ . Doubling the one tail, the p-value is 0.04. (Students may

have approximate results, and a small number of students may have a p-value of about 0.05.) Since the p-value is low, we reject  $H_0$ . The data provide strong evidence that people react differently under the two scenarios.

## 7 Introduction to linear regression

**7.1** (a) The residual plot will show randomly distributed residuals around 0. The variance is also approximately constant. (b) The residuals will show a fan shape, with higher variability for smaller  $x$ . There will also be many points on the right above the line. There is trouble with the model being fit here.

**7.3** (a) Strong relationship, but a straight line would not fit the data. (b) Strong relationship, and a linear fit would be reasonable. (c) Weak relationship, and trying a linear fit would be reasonable. (d) Moderate relationship, but a straight line would not fit the data. (e) Strong relationship, and a linear fit would be reasonable. (f) Weak relationship, and trying a linear fit would be reasonable.

**7.5** (a) Exam 2 since there is less of a scatter in the plot of final exam grade versus exam 2. Notice that the relationship between Exam 1 and the Final Exam appears to be slightly nonlinear. (b) Exam 2 and the final are relatively close to each other chronologically, or Exam 2 may be cumulative so has greater similarities in material to the final exam. Answers may vary for part (b).

**7.7** (a)  $r = -0.7 \rightarrow (4)$ . (b)  $r = 0.45 \rightarrow (3)$ . (c)  $r = 0.06 \rightarrow (1)$ . (d)  $r = 0.92 \rightarrow (2)$ .

**7.9** (a) True. (b) False, correlation is a measure of the linear association between any two numerical variables.

**7.11** (a) The relationship is positive, weak, and possibly linear. However, there do appear to be some anomalous observations along the left where several students have the same height that is notably far from the cloud of the other points. Additionally, there are many students who appear not to have driven a car, and they are represented by a set of points along the bottom of the scatterplot. (b) There is no obvious explanation why simply being tall should lead a

person to drive faster. However, one confounding factor is gender. Males tend to be taller than females on average, and personal experiences (anecdotal) may suggest they drive faster. If we were to follow-up on this suspicion, we would find that sociological studies confirm this suspicion. (c) Males are taller on average and they drive faster. The gender variable is indeed an important confounding variable.

**7.13** (a) There is a somewhat weak, positive, possibly linear relationship between the distance traveled and travel time. There is clustering near the lower left corner that we should take special note of. (b) Changing the units will not change the form, direction or strength of the relationship between the two variables. If longer distances measured in miles are associated with longer travel time measured in minutes, longer distances measured in kilometers will be associated with longer travel time measured in hours. (c) Changing units doesn't affect correlation:  $r = 0.636$ .

**7.15** (a) There is a moderate, positive, and linear relationship between shoulder girth and height. (b) Changing the units, even if just for one of the variables, will not change the form, direction or strength of the relationship between the two variables.

**7.17** In each part, we can write the husband ages as a linear function of the wife ages.

- (a)  $age_H = age_W + 3$ .
- (b)  $age_H = age_W - 2$ .
- (c)  $age_H = 2 \times age_W$ .

Since the slopes are positive and these are perfect linear relationships, the correlation will be exactly 1 in all three parts. An alternative way to gain insight into this solution is to create a mock data set, e.g. 5 women aged 26, 27, 28, 29, and 30, then find the husband ages for each wife in each part and create a scatterplot.

**7.19** Correlation: no units. Intercept: kg.  
Slope: kg/cm.

**7.21** Over-estimate. Since the residual is calculated as *observed* – *predicted*, a negative residual means that the predicted value is higher than the observed value.

**7.23** (a) There is a positive, very strong, linear association between the number of tourists and spending. (b) Explanatory: number of tourists (in thousands). Response: spending (in millions of US dollars). (c) We can predict spending for a given number of tourists using a regression line. This may be useful information for determining how much the country may want to spend in advertising abroad, or to forecast expected revenues from tourism. (d) Even though the relationship appears linear in the scatterplot, the residual plot actually shows a nonlinear relationship. This is not a contradiction: residual plots can show divergences from linearity that can be difficult to see in a scatterplot. A simple linear model is inadequate for modeling these data. It is also important to consider that these data are observed sequentially, which means there may be a hidden structure not evident in the current plots but that is important to consider.

**7.25** (a) First calculate the slope:  $b_1 = R \times s_y/s_x = 0.636 \times 113/99 = 0.726$ . Next, make use of the fact that the regression line passes through the point  $(\bar{x}, \bar{y})$ :  $\bar{y} = b_0 + b_1 \times \bar{x}$ . Plug in  $\bar{x}$ ,  $\bar{y}$ , and  $b_1$ , and solve for  $b_0$ : 51. Solution: *travel time* = 51 + 0.726 × *distance*. (b)  $b_1$ : For each additional mile in distance, the model predicts an additional 0.726 minutes in travel time.  $b_0$ : When the distance traveled is 0 miles, the travel time is expected to be 51 minutes. It does not make sense to have a travel distance of 0 miles in this context. Here, the *y*-intercept serves only to adjust the height of the line and is meaningless by itself. (c)  $R^2 = 0.636^2 = 0.40$ . About 40% of the variability in travel time is accounted for by the model, i.e. explained by the distance traveled. (d)  $\widehat{\text{travel time}} = 51 + 0.726 \times \text{distance} = 51 + 0.726 \times 103 \approx 126$  minutes. (Note: we should be cautious in our predictions with this model since we have not yet evaluated whether it is a well-fit model.) (e)  $e_i = y_i - \hat{y}_i = 168 - 126 = 42$  minutes. A

positive residual means that the model underestimates the travel time. (f) No, this calculation would require extrapolation.

**7.27** There is an upwards trend. However, the variability is higher for higher calorie counts, and it looks like there might be two clusters of observations above and below the line on the right, so we should be cautious about fitting a linear model to these data.

**7.29** (a)  $\widehat{\text{murder}} = -29.901 + 2.559 \times \text{poverty}\%$   
(b) Expected murder rate in metropolitan areas with no poverty is -29.901 per million. This is obviously not a meaningful value, it just serves to adjust the height of the regression line.  
(c) For each additional percentage increase in poverty, we expect murders per million to be higher on average by 2.559. (d) Poverty level explains 70.52% of the variability in murder rates in metropolitan areas. (e)  $\sqrt{0.7052} = 0.8398$

**7.31** (a) There is an outlier in the bottom right. Since it is far from the center of the data, it is a point with high leverage. It is also an influential point since, without that observation, the regression line would have a very different slope.

(b) There is an outlier in the bottom right. Since it is far from the center of the data, it is a point with high leverage. However, it does not appear to be affecting the line much, so it is not an influential point.

(c) The observation is in the center of the data (in the x-axis direction), so this point does *not* have high leverage. This means the point won't have much effect on the slope of the line and so is not an influential point.

**7.33** (a) There is a negative, moderate-to-strong, somewhat linear relationship between percent of families who own their home and the percent of the population living in urban areas in 2010. There is one outlier: a state where 100% of the population is urban. The variability in the percent of homeownership also increases as we move from left to right in the plot. (b) The outlier is located in the bottom right corner, horizontally far from the center of the other points, so it is a point with high leverage. It is an influential point since excluding this point from the analysis would greatly affect the slope of the regression line.

**7.35** (a) The relationship is positive, moderate-to-strong, and linear. There are a few outliers but no points that appear to be influential. (b)  $\widehat{weight} = -105.0113 + 1.0176 \times height$ . Slope: For each additional centimeter in height, the model predicts the average weight to be 1.0176 additional kilograms (about 2.2 pounds). Intercept: People who are 0 centimeters tall are expected to weigh -105.0113 kilograms. This is obviously not possible. Here, the  $y$ -intercept serves only to adjust the height of the line and is meaningless by itself. (c)  $H_0$ : The true slope coefficient of height is zero ( $\beta_1 = 0$ ).  $H_A$ : The true slope coefficient of height is greater than zero ( $\beta_1 > 0$ ). A two-sided test would also be acceptable for this application. The p-value for the two-sided alternative hypothesis ( $\beta_1 \neq 0$ ) is incredibly small, so the p-value for the one-sided hypothesis will be even smaller. That is, we reject  $H_0$ . The data provide convincing evidence that height and weight are positively correlated. The true slope parameter is indeed greater than 0. (d)  $R^2 = 0.72^2 = 0.52$ . Approximately 52% of the variability in weight can be explained by the height of individuals.

**7.37** (a)  $H_0: \beta_1 = 0$ .  $H_A: \beta_1 > 0$ . A two-sided test would also be acceptable for this application. The p-value, as reported in the table, is incredibly small. Thus, for a one-sided test, the p-value will also be incredibly small, and we reject  $H_0$ . The data provide convincing evidence that wives' and husbands' heights are positively correlated. (b)  $height_W = 43.5755 + 0.2863 \times height_H$ . (c) Slope: For each additional inch

in husband's height, the average wife's height is expected to be an additional 0.2863 inches on average. Intercept: Men who are 0 inches tall are expected to have wives who are, on average, 43.5755 inches tall. The intercept here is meaningless, and it serves only to adjust the height of the line. (d) The slope is positive, so  $r$  must also be positive.  $r = \sqrt{0.09} = 0.30$ . (e) 63.2612. Since  $R^2$  is low, the prediction based on this regression model is not very reliable. (f) No, we should avoid extrapolating.

**7.39** (a)  $r = \sqrt{0.28} \approx -0.53$ . We know the correlation is negative due to the negative association shown in the scatterplot. (b) The residuals appear to be fan shaped, indicating non-constant variance. Therefore a simple least squares fit is not appropriate for these data.

**7.41** (a)  $H_0 : \beta_1 = 0$ ;  $H_A : \beta_1 \neq 0$  (b) The p-value for this test is approximately 0, therefore we reject  $H_0$ . The data provide convincing evidence that poverty percentage is a significant predictor of murder rate. (c)  $n = 20, df = 18, T_{18}^* = 2.10; 2.559 \pm 2.10 \times 0.390 = (1.74, 3.378)$ ; For each percentage point poverty is higher, murder rate is expected to be higher on average by 1.74 to 3.378 per million. (d) Yes, we rejected  $H_0$  and the confidence interval does not include 0.

**7.43** This is a one-sided test, so the p-value should be half of the p-value given in the regression table, which will be approximately 0. Therefore the data provide convincing evidence that poverty percentage is positively associated with murder rate.

## 8 Multiple and logistic regression

**??** (a)  $\widehat{baby\_weight} = 123.05 - 8.94 \times smoke$  (b) The estimated body weight of babies born to smoking mothers is 8.94 ounces lower than babies born to non-smoking mothers. Smoker:  $123.05 - 8.94 \times 1 = 114.11$  ounces. Non-smoker:  $123.05 - 8.94 \times 0 = 123.05$  ounces. (c)  $H_0: \beta_1 = 0$ .  $H_A: \beta_1 \neq 0$ .  $T = -8.65$ , and the p-value is approximately 0. Since the p-value is very small, we reject  $H_0$ . The data provide strong evidence that the true slope parameter is different than 0 and that there is an association between birth weight and smoking. Furthermore, having rejected  $H_0$ , we can conclude that smoking is associated with lower birth weights.

**??** (a)  $\widehat{baby\_weight} = -80.41 + 0.44 \times gestation - 3.33 \times parity - 0.01 \times age + 1.15 \times height + 0.05 \times weight - 8.40 \times smoke$ . (b)  $\beta_{gestation}$ : The model predicts a 0.44 ounce increase in the birth weight of the baby for each additional day of pregnancy, all else held constant.  $\beta_{age}$ : The model predicts a 0.01 ounce decrease in the birth weight of the baby for each additional year in mother's age, all else held constant. (c) Parity might be correlated with one of the other variables in the model, which complicates model estimation. (d)  $\widehat{baby\_weight} = 120.58$ .  $e = 120 - 120.58 = -0.58$ . The model over-predicts this baby's birth weight. (e)  $R^2 = 0.2504$ .  $R^2_{adj} = 0.2468$ .

?? (a) (-0.32, 0.16). We are 95% confident that male students on average have GPAs 0.32 points lower to 0.16 points higher than females when controlling for the other variables in the model.  
 (b) Yes, since the p-value is larger than 0.05 in all cases (not including the intercept).

?? Remove age.

?? Based on the p-value alone, either gestation or smoke should be added to the model first. However, since the adjusted  $R^2$  for the model with gestation is higher, it would be preferable to add gestation in the first step of the forward-selection algorithm. (Other explanations are possible. For instance, it would be reasonable to only use the adjusted  $R^2$ .)

?? She should use p-value selection since she is interested in finding out about significant predictors, not just optimizing predictions.

?? Nearly normal residuals: The normal probability plot shows a nearly normal distribution of the residuals, however, there are some minor irregularities at the tails. With a data set so large, these would not be a concern.

Constant variability of residuals: The scatterplot of the residuals versus the fitted values does not show any overall structure. However, values that have very low or very high fitted values appear to also have somewhat larger outliers. In addition, the residuals do appear to have constant variability between the two parity and smoking status groups, though these items are relatively minor.

Independent residuals: The scatterplot of residuals versus the order of data collection shows a random scatter, suggesting that there is no apparent structures related to the order the data were collected.

Linear relationships between the response variable and numerical explanatory variables: The residuals vs. height and weight of mother are randomly distributed around 0. The residuals

vs. length of gestation plot also does not show any clear or strong remaining structures, with the possible exception of very short or long gestations. The rest of the residuals do appear to be randomly distributed around 0.

All concerns raised here are relatively mild. There are some outliers, but there is so much data that the influence of such observations will be minor.

?? (a) There are a few potential outliers, e.g. on the left in the `total_length` variable, but nothing that will be of serious concern in a data set this large. (b) When coefficient estimates are sensitive to which variables are included in the model, this typically indicates that some variables are collinear. For example, a possum's gender may be related to its head length, which would explain why the coefficient (and p-value) for `sex_male` changed when we removed the `head_length` variable. Likewise, a possum's skull width is likely to be related to its head length, probably even much more closely related than the head length was to gender.

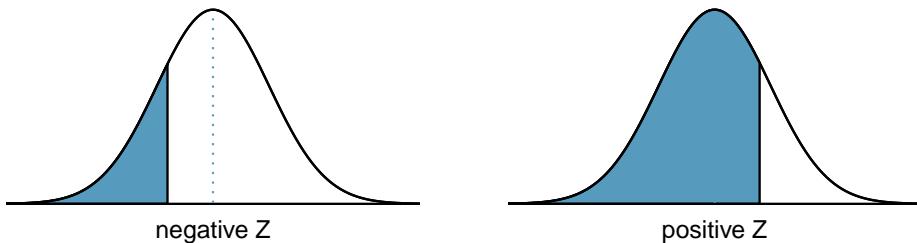
?? (a) The logistic model relating  $\hat{p}_i$  to the predictors may be written as  $\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) = 33.5095 - 1.4207 \times \text{sex\_male}_i - 0.2787 \times \text{skull\_width}_i + 0.5687 \times \text{total\_length}_i - 1.8057 \times \text{tail\_length}_i$ . Only `total_length` has a positive association with a possum being from Victoria.  
 (b)  $\hat{p} = 0.0062$ . While the probability is very near zero, we have not run diagnostics on the model. We might also be a little skeptical that the model will remain accurate for a possum found in a US zoo. For example, perhaps the zoo selected a possum with specific characteristics but only looked in one region. On the other hand, it is encouraging that the possum was caught in the wild. (Answers regarding the reliability of the model probability will vary.)

# Appendix B

## Distribution tables

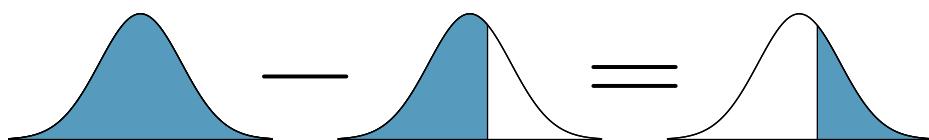
### B.1 Normal Probability Table

The area to the left of  $Z$  represents the percentile of the observation. The normal probability table always lists percentiles.

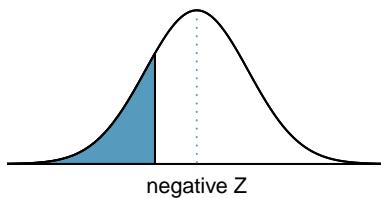


To find the area to the right, calculate 1 minus the area to the left.

$$1.0000 - 0.6664 = 0.3336$$

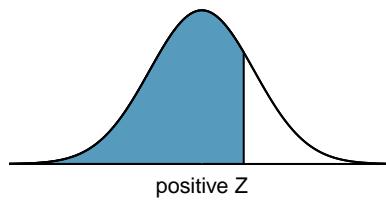


For additional details about working with the normal distribution and the normal probability table, see Section 3.1, which starts on page 76.



Second decimal place of $Z$										$Z$
0.09	0.08	0.07	0.06	0.05	0.04	0.03	0.02	0.01	0.00	
0.0002	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	-3.4
0.0003	0.0004	0.0004	0.0004	0.0004	0.0004	0.0004	0.0005	0.0005	0.0005	-3.3
0.0005	0.0005	0.0005	0.0006	0.0006	0.0006	0.0006	0.0006	0.0007	0.0007	-3.2
0.0007	0.0007	0.0008	0.0008	0.0008	0.0008	0.0009	0.0009	0.0009	0.0010	-3.1
0.0010	0.0010	0.0011	0.0011	0.0011	0.0012	0.0012	0.0013	0.0013	0.0013	-3.0
0.0014	0.0014	0.0015	0.0015	0.0016	0.0016	0.0017	0.0018	0.0018	0.0019	-2.9
0.0019	0.0020	0.0021	0.0021	0.0022	0.0023	0.0023	0.0024	0.0025	0.0026	-2.8
0.0026	0.0027	0.0028	0.0029	0.0030	0.0031	0.0032	0.0033	0.0034	0.0035	-2.7
0.0036	0.0037	0.0038	0.0039	0.0040	0.0041	0.0043	0.0044	0.0045	0.0047	-2.6
0.0048	0.0049	0.0051	0.0052	0.0054	0.0055	0.0057	0.0059	0.0060	0.0062	-2.5
0.0064	0.0066	0.0068	0.0069	0.0071	0.0073	0.0075	0.0078	0.0080	0.0082	-2.4
0.0084	0.0087	0.0089	0.0091	0.0094	0.0096	0.0099	0.0102	0.0104	0.0107	-2.3
0.0110	0.0113	0.0116	0.0119	0.0122	0.0125	0.0129	0.0132	0.0136	0.0139	-2.2
0.0143	0.0146	0.0150	0.0154	0.0158	0.0162	0.0166	0.0170	0.0174	0.0179	-2.1
0.0183	0.0188	0.0192	0.0197	0.0202	0.0207	0.0212	0.0217	0.0222	0.0228	-2.0
0.0233	0.0239	0.0244	0.0250	0.0256	0.0262	0.0268	0.0274	0.0281	0.0287	-1.9
0.0294	0.0301	0.0307	0.0314	0.0322	0.0329	0.0336	0.0344	0.0351	0.0359	-1.8
0.0367	0.0375	0.0384	0.0392	0.0401	0.0409	0.0418	0.0427	0.0436	0.0446	-1.7
0.0455	0.0465	0.0475	0.0485	0.0495	0.0505	0.0516	0.0526	0.0537	0.0548	-1.6
0.0559	0.0571	0.0582	0.0594	0.0606	0.0618	0.0630	0.0643	0.0655	0.0668	-1.5
0.0681	0.0694	0.0708	0.0721	0.0735	0.0749	0.0764	0.0778	0.0793	0.0808	-1.4
0.0823	0.0838	0.0853	0.0869	0.0885	0.0901	0.0918	0.0934	0.0951	0.0968	-1.3
0.0985	0.1003	0.1020	0.1038	0.1056	0.1075	0.1093	0.1112	0.1131	0.1151	-1.2
0.1170	0.1190	0.1210	0.1230	0.1251	0.1271	0.1292	0.1314	0.1335	0.1357	-1.1
0.1379	0.1401	0.1423	0.1446	0.1469	0.1492	0.1515	0.1539	0.1562	0.1587	-1.0
0.1611	0.1635	0.1660	0.1685	0.1711	0.1736	0.1762	0.1788	0.1814	0.1841	-0.9
0.1867	0.1894	0.1922	0.1949	0.1977	0.2005	0.2033	0.2061	0.2090	0.2119	-0.8
0.2148	0.2177	0.2206	0.2236	0.2266	0.2296	0.2327	0.2358	0.2389	0.2420	-0.7
0.2451	0.2483	0.2514	0.2546	0.2578	0.2611	0.2643	0.2676	0.2709	0.2743	-0.6
0.2776	0.2810	0.2843	0.2877	0.2912	0.2946	0.2981	0.3015	0.3050	0.3085	-0.5
0.3121	0.3156	0.3192	0.3228	0.3264	0.3300	0.3336	0.3372	0.3409	0.3446	-0.4
0.3483	0.3520	0.3557	0.3594	0.3632	0.3669	0.3707	0.3745	0.3783	0.3821	-0.3
0.3859	0.3897	0.3936	0.3974	0.4013	0.4052	0.4090	0.4129	0.4168	0.4207	-0.2
0.4247	0.4286	0.4325	0.4364	0.4404	0.4443	0.4483	0.4522	0.4562	0.4602	-0.1
0.4641	0.4681	0.4721	0.4761	0.4801	0.4840	0.4880	0.4920	0.4960	0.5000	-0.0

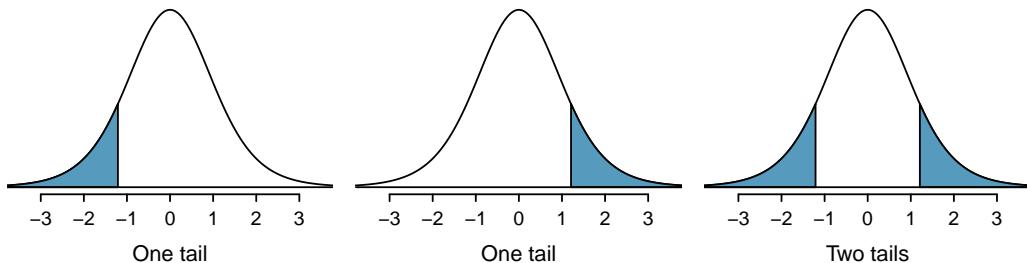
\*For  $Z \leq -3.50$ , the probability is less than or equal to 0.0002.



Z	Second decimal place of Z									
	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9279	0.9292	0.9306	0.9319
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441
1.6	0.9452	0.9463	0.9474	0.9484	0.9495	0.9505	0.9515	0.9525	0.9535	0.9545
1.7	0.9554	0.9564	0.9573	0.9582	0.9591	0.9599	0.9608	0.9616	0.9625	0.9633
1.8	0.9641	0.9649	0.9656	0.9664	0.9671	0.9678	0.9686	0.9693	0.9699	0.9706
1.9	0.9713	0.9719	0.9726	0.9732	0.9738	0.9744	0.9750	0.9756	0.9761	0.9767
2.0	0.9772	0.9778	0.9783	0.9788	0.9793	0.9798	0.9803	0.9808	0.9812	0.9817
2.1	0.9821	0.9826	0.9830	0.9834	0.9838	0.9842	0.9846	0.9850	0.9854	0.9857
2.2	0.9861	0.9864	0.9868	0.9871	0.9875	0.9878	0.9881	0.9884	0.9887	0.9890
2.3	0.9893	0.9896	0.9898	0.9901	0.9904	0.9906	0.9909	0.9911	0.9913	0.9916
2.4	0.9918	0.9920	0.9922	0.9925	0.9927	0.9929	0.9931	0.9932	0.9934	0.9936
2.5	0.9938	0.9940	0.9941	0.9943	0.9945	0.9946	0.9948	0.9949	0.9951	0.9952
2.6	0.9953	0.9955	0.9956	0.9957	0.9959	0.9960	0.9961	0.9962	0.9963	0.9964
2.7	0.9965	0.9966	0.9967	0.9968	0.9969	0.9970	0.9971	0.9972	0.9973	0.9974
2.8	0.9974	0.9975	0.9976	0.9977	0.9977	0.9978	0.9979	0.9979	0.9980	0.9981
2.9	0.9981	0.9982	0.9982	0.9983	0.9984	0.9984	0.9985	0.9985	0.9986	0.9986
3.0	0.9987	0.9987	0.9987	0.9988	0.9988	0.9989	0.9989	0.9989	0.9990	0.9990
3.1	0.9990	0.9991	0.9991	0.9991	0.9992	0.9992	0.9992	0.9992	0.9993	0.9993
3.2	0.9993	0.9993	0.9994	0.9994	0.9994	0.9994	0.9994	0.9995	0.9995	0.9995
3.3	0.9995	0.9995	0.9995	0.9996	0.9996	0.9996	0.9996	0.9996	0.9996	0.9997
3.4	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9998

\*For  $Z \geq 3.50$ , the probability is greater than or equal to 0.9998.

## B.2 t-Probability Table

Figure B.1: Tails for the  $t$ -distribution.

	one tail	0.100	0.050	0.025	0.010	0.005
	two tails	0.200	0.100	0.050	0.020	0.010
df	1	3.08	6.31	12.71	31.82	63.66
	2	1.89	2.92	4.30	6.96	9.92
	3	1.64	2.35	3.18	4.54	5.84
	4	1.53	2.13	2.78	3.75	4.60
	5	1.48	2.02	2.57	3.36	4.03
	6	1.44	1.94	2.45	3.14	3.71
	7	1.41	1.89	2.36	3.00	3.50
	8	1.40	1.86	2.31	2.90	3.36
	9	1.38	1.83	2.26	2.82	3.25
	10	1.37	1.81	2.23	2.76	3.17
	11	1.36	1.80	2.20	2.72	3.11
	12	1.36	1.78	2.18	2.68	3.05
	13	1.35	1.77	2.16	2.65	3.01
	14	1.35	1.76	2.14	2.62	2.98
	15	1.34	1.75	2.13	2.60	2.95
	16	1.34	1.75	2.12	2.58	2.92
	17	1.33	1.74	2.11	2.57	2.90
	18	1.33	1.73	2.10	2.55	2.88
	19	1.33	1.73	2.09	2.54	2.86
	20	1.33	1.72	2.09	2.53	2.85
	21	1.32	1.72	2.08	2.52	2.83
	22	1.32	1.72	2.07	2.51	2.82
	23	1.32	1.71	2.07	2.50	2.81
	24	1.32	1.71	2.06	2.49	2.80
	25	1.32	1.71	2.06	2.49	2.79
	26	1.31	1.71	2.06	2.48	2.78
	27	1.31	1.70	2.05	2.47	2.77
	28	1.31	1.70	2.05	2.47	2.76
	29	1.31	1.70	2.05	2.46	2.76
	30	1.31	1.70	2.04	2.46	2.75

	one tail	0.100	0.050	0.025	0.010	0.005
	two tails	0.200	0.100	0.050	0.020	0.010
df	31	1.31	1.70	2.04	2.45	2.74
	32	1.31	1.69	2.04	2.45	2.74
	33	1.31	1.69	2.03	2.44	2.73
	34	1.31	1.69	2.03	2.44	2.73
	35	1.31	1.69	2.03	2.44	2.72
	36	1.31	1.69	2.03	2.43	2.72
	37	1.30	1.69	2.03	2.43	2.72
	38	1.30	1.69	2.02	2.43	2.71
	39	1.30	1.68	2.02	2.43	2.71
	40	1.30	1.68	2.02	2.42	2.70
	41	1.30	1.68	2.02	2.42	2.70
	42	1.30	1.68	2.02	2.42	2.70
	43	1.30	1.68	2.02	2.42	2.70
	44	1.30	1.68	2.02	2.41	2.69
	45	1.30	1.68	2.01	2.41	2.69
	46	1.30	1.68	2.01	2.41	2.69
	47	1.30	1.68	2.01	2.41	2.68
	48	1.30	1.68	2.01	2.41	2.68
	49	1.30	1.68	2.01	2.40	2.68
	50	1.30	1.68	2.01	2.40	2.68
	60	1.30	1.67	2.00	2.39	2.66
	70	1.29	1.67	1.99	2.38	2.65
	80	1.29	1.66	1.99	2.37	2.64
	90	1.29	1.66	1.99	2.37	2.63
	100	1.29	1.66	1.98	2.36	2.63
	150	1.29	1.66	1.98	2.35	2.61
	200	1.29	1.65	1.97	2.35	2.60
	300	1.28	1.65	1.97	2.34	2.59
	400	1.28	1.65	1.97	2.34	2.59
	500	1.28	1.65	1.96	2.33	2.59
	$\infty$	1.28	1.65	1.96	2.33	2.58

### B.3 Chi-Square Probability Table

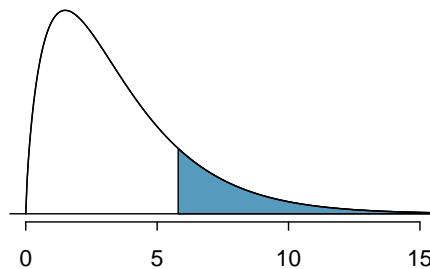


Figure B.2: Areas in the chi-square table always refer to the right tail.

Upper tail	0.3	0.2	0.1	0.05	0.02	0.01	0.005	0.001
df	1	2	3	4	5	6	7	8
1	1.07	1.64	2.71	3.84	5.41	6.63	7.88	10.83
2	2.41	3.22	4.61	5.99	7.82	9.21	10.60	13.82
3	3.66	4.64	6.25	7.81	9.84	11.34	12.84	16.27
4	4.88	5.99	7.78	9.49	11.67	13.28	14.86	18.47
5	6.06	7.29	9.24	11.07	13.39	15.09	16.75	20.52
6	7.23	8.56	10.64	12.59	15.03	16.81	18.55	22.46
7	8.38	9.80	12.02	14.07	16.62	18.48	20.28	24.32
8	9.52	11.03	13.36	15.51	18.17	20.09	21.95	26.12
9	10.66	12.24	14.68	16.92	19.68	21.67	23.59	27.88
10	11.78	13.44	15.99	18.31	21.16	23.21	25.19	29.59
11	12.90	14.63	17.28	19.68	22.62	24.72	26.76	31.26
12	14.01	15.81	18.55	21.03	24.05	26.22	28.30	32.91
13	15.12	16.98	19.81	22.36	25.47	27.69	29.82	34.53
14	16.22	18.15	21.06	23.68	26.87	29.14	31.32	36.12
15	17.32	19.31	22.31	25.00	28.26	30.58	32.80	37.70
16	18.42	20.47	23.54	26.30	29.63	32.00	34.27	39.25
17	19.51	21.61	24.77	27.59	31.00	33.41	35.72	40.79
18	20.60	22.76	25.99	28.87	32.35	34.81	37.16	42.31
19	21.69	23.90	27.20	30.14	33.69	36.19	38.58	43.82
20	22.77	25.04	28.41	31.41	35.02	37.57	40.00	45.31
25	28.17	30.68	34.38	37.65	41.57	44.31	46.93	52.62
30	33.53	36.25	40.26	43.77	47.96	50.89	53.67	59.70
40	44.16	47.27	51.81	55.76	60.44	63.69	66.77	73.40
50	54.72	58.16	63.17	67.50	72.61	76.15	79.49	86.66