

CONTENTS

<i>Preface</i>		vii
1	Introduction	1
2	Working with Data	9
3	Summarizing Data	17
4	Data Visualization	28
5	Fitting Models to Data	40
6	Probability	61
7	Sampling	80
8	Resampling and Simulation	86
9	Hypothesis Testing	93
10	Quantifying Effects and Designing Studies	111
11	Bayesian Statistics	124
12	Modeling Categorical Relationships	138
13	Modeling Continuous Relationships	147
14	The General Linear Model	158
15	Comparing Means	178
16	Multivariate Statistics	190

17	Practical Statistical Modeling	211
18	Doing Reproducible Research	248
	<i>Bibliography</i>	259
	<i>Index</i>	263

1

Introduction

Learning Objectives

Having read this chapter, you should be able to

- Describe the central goals and fundamental concepts of statistics.
- Describe the difference between experimental and observational research with regard to what can be inferred about causality.
- Explain how randomization provides the ability to make inferences about causation.

What Is Statistical Thinking?

Statistical thinking is a way of understanding a complex world by describing it in relatively simple terms that nonetheless capture essential aspects of its structure or function, and that also provide us with some idea of how uncertain we are about that knowledge. The foundations of statistical thinking come primarily from mathematics and statistics but also from computer science, psychology, and other fields of study.

We can distinguish statistical thinking from other forms of thinking that are less likely to describe the world accurately. In particular, human intuition frequently tries to answer the same questions that we can answer using statistical thinking, but it often gets the answer wrong. For example, in recent years most Americans have reported that they think violent crime is worse in the current year compared to the previous year (Pew 2020). However, a statistical analysis of the actual crime data showed that in fact violent crime was steadily *decreasing* during that time. Intuition fails us because we rely on best guesses (which psychologists refer to as *heuristics*) that can often get it wrong. For example, humans often judge the prevalence of some event (like violent crime) using an *availability heuristic*—that is, how easily we can think of an example of violent crime. For this reason, our judgments of increasing crime rates may be more reflective of increasing news coverage, in spite of an actual decrease in the rate of crime. Statistical thinking provides us with the tools to more accurately understand the world and overcome the biases of human judgment.

Dealing with Statistics Anxiety

Many people come to their first statistics class with a lot of trepidation and anxiety, especially once they hear that they will also have to learn to code in order to analyze data. In my class, I give students a survey prior to the first session to measure their attitude toward statistics, asking them to rate a number of statements on a scale of 1 (strongly disagree) to 7 (strongly agree). One of the items on the survey is “The thought of being enrolled in a statistics course makes me nervous.” In a recent class, almost two-thirds of the students responded with a 5 or higher, and about one-fourth of the students said they strongly agreed with the statement. So if you feel nervous about starting to learn statistics, you are not alone.

Anxiety feels uncomfortable, but psychology tells us that this kind of emotional arousal can actually help us perform *better* on many tasks, by focusing our attention. So if you start to feel anxious about the material in this book, remind yourself that many other readers are feeling similarly, and that this emotional arousal could actually help cement the material in your brain more effectively (even if it doesn’t seem like it!).

What Can Statistics Do for Us?

There are three major things that we can do with statistics:

- *Describe*: The world is complex and we often need to describe it in a simplified way that we can understand.
- *Decide*: We often need to make decisions based on data, usually in the face of uncertainty.
- *Predict*: We often wish to make predictions about new situations based on our knowledge of previous situations.

Let’s look at an example of these in action, centered on a question that many of us are interested in: How do we decide what’s healthy to eat? There are many different sources of guidance: government dietary guidelines, diet books, and bloggers, just to name a few. Let’s focus in on a specific question: Is saturated fat in our diet a bad thing?

One way that we might answer this question is common sense. If we eat fat, then it’s going to turn straight into fat in our bodies, right? And we have all seen photos of arteries clogged with fat, so eating fat is going to clog our arteries, right?

Another way that we might answer this question is by listening to authority figures. The dietary guidelines from the US Food and Drug Administration have as one of their key recommendations that “A healthy eating pattern limits saturated fats.” You might hope that these guidelines would be based on good science, and in some cases they are; but as Nina Teicholz (2014) outlined in her book *Big Fat Surprise*, this particular recommendation seems to be based more on the long-standing dogma of nutrition researchers than on actual evidence.

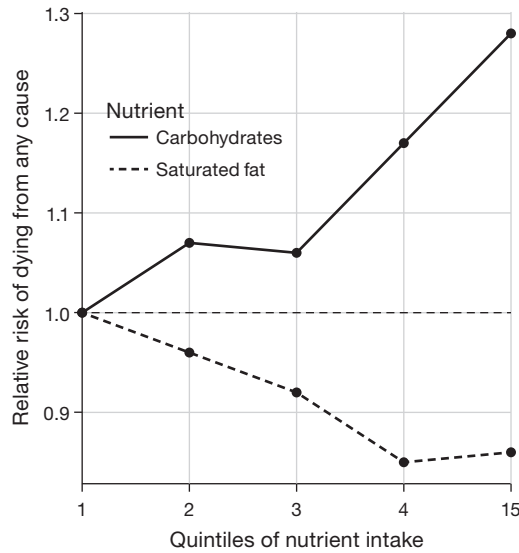


FIGURE 1.1. A plot of data from the PURE study, showing the relationship between death from any cause and the relative intake of saturated fats and carbohydrates.

Finally, we might look at actual scientific research. Let's start by looking at the PURE study, which has examined diets and health outcomes (including death) in more than 135,000 people from 18 different countries. In one of the analyses of this dataset (Dehghan et al. 2017, published in *the Lancet*), the PURE investigators reported an analysis of how intake of various classes of macronutrients (including saturated fats and carbohydrates) was related to the likelihood of people dying during the time they participated in the study. Participants were followed for a *median* of 7.4 years, meaning that half of the people in the study were followed for less than 7.4 years and half were followed for more than 7.4 years. Figure 1.1 plots some of the data from the study (extracted from the publication), showing the relationship between the intake of both saturated fats and carbohydrates and the risk of dying from any cause.

This plot is based on 10 numbers. To obtain these numbers, the researchers split the group of 135,335 study participants (which we call the *sample*) into five groups (quintiles) after ordering them in terms of their intake of either of the nutrients; the first quintile contains the 20% of people with the lowest intake, and the fifth quintile contains the 20% with the highest intake. The researchers then computed how often people in each of those groups died during the time they were being followed. The figure expresses this in terms of the *relative risk* of dying in comparison to the lowest quintile: if this number is greater than one, it means that people in the group are *more* likely to die than are people in the lowest quintile, whereas if it's less than one, it means that people in the group are *less* likely to die. The figure is pretty clear: people who ate more saturated fat were *less* likely to die during the study, with the lowest death rate seen for people who were in the fourth quintile (that

is, those who ate more fat than the lowest 60% but less than the top 20%). The opposite is seen for carbohydrates; the more carbs a person ate, the more likely they were to die during the study. This example shows how we can use statistics to *describe* a complex dataset in terms of a much simpler set of numbers; if we had to look at the data from each of the study participants at the same time, we would be overloaded with data and it would be hard to see the pattern that emerges when they are described more simply.

The numbers in figure 1.1 seem to show that deaths decrease with saturated fat and increase with carbohydrate intake, but we also know that there is a lot of uncertainty in the data; there are some people who died early even though they ate a low-carb diet, and, similarly, some people who ate a ton of carbs but lived to a ripe old age. Given this variability, we want to *decide* whether the relationships we see in the data are large enough that we wouldn't expect them to occur randomly if there was not truly a relationship between diet and longevity. Statistics provide us with the tools to make these kinds of decisions, and often people from the outside view this as *the* main purpose of statistics. But as we will see throughout the book, this need for black-and-white decisions based on fuzzy evidence has often led researchers astray.

Based on the data, we would also like to be able to *predict* future outcomes. For example, a life insurance company might want to use data about a particular person's intake of fat and carbohydrates to predict how long they are likely to live. An important aspect of prediction is that it requires us to generalize from the data we already have to some other situation, often in the future; if our conclusions were limited to the specific people in the study at a particular time, then the study would not be very useful. In general, researchers must assume that their particular sample is representative of a larger *population*, which requires that they obtain the sample in a way that provides an unbiased picture of the population. For example, if the PURE study had recruited all of its participants from religious sects that practice vegetarianism, then we probably wouldn't want to generalize the results to people who follow different dietary standards.

The Big Ideas of Statistics

There are a number of very basic ideas that cut through nearly all aspects of statistical thinking. Several of these are outlined by Stigler (2016) in his outstanding book *The Seven Pillars of Statistical Wisdom*, which I have augmented here.

Learning from Data

One way to think of statistics is as a set of tools that enable us to learn from data. In any situation, we start with a set of ideas or *hypotheses* about what might be the case. In the PURE study, the researchers may have started out with the expectation that eating more fat would lead to higher death rates, given the prevailing negative dogma about saturated fats. Later in the book we introduce the idea of *prior knowledge*, that is, the knowledge

that we bring to a situation. This prior knowledge can vary in its strength, often based on our level of experience; if I visit a restaurant for the first time, I am likely to have a weak expectation of how good it will be, but if I visit a restaurant where I have eaten 10 times before, my expectations will be much stronger. Similarly, if I look at a restaurant review site and see that a restaurant's average rating of four stars is based on only three reviews, I will have a weaker expectation than I would if it were based on three hundred reviews.

Statistics provides us with a way to describe how new data can be best used to update our beliefs, and in this way there are deep links between statistics and psychology. In fact, many theories of human and animal learning from psychology are closely aligned with ideas from the field of *machine learning*—a new field at the interface of statistics and computer science that focuses on how to build computer algorithms that can learn from experience. While statistics and machine learning often try to solve the same problems, researchers from these fields frequently take very different approaches; the famous statistician Leo Breiman (2001) once referred to them as “the two cultures” to reflect how different their approaches can be. In this book, I try to blend the two cultures together because both approaches provide useful tools for thinking about data.

Aggregation

Another way to think of statistics is as “the science of throwing away data.” In the example of the PURE study above, we took more than 100,000 numbers and condensed them into 10. It is this kind of *aggregation* that is one of the most important concepts in statistics. When it was first advanced, this idea was revolutionary: if we throw out all the details about every one of the participants, then how can we be sure we aren't missing something important?

As we will see, statistics provides us with ways to characterize the structure of aggregations of data, with theoretical foundations that explain why this usually works well. However, it's also important to keep in mind that aggregation can go too far, and later we will encounter cases where a summary can provide a very misleading picture of the data being summarized.

Uncertainty

The world is an uncertain place. We now know that cigarette smoking causes lung cancer, but this causation is probabilistic: a 68-year-old man who smoked two packs a day for the past 50 years and continues to smoke has a 15% (1 out of 7) risk of getting lung cancer, which is much higher than the chance of lung cancer in a nonsmoker. However, it also means that there will be many people who smoke their entire lives and never get lung cancer. Statistics provides us with the tools to characterize uncertainty, to make decisions under uncertainty, and to make predictions whose uncertainty we can quantify.

One often sees journalists write that scientific researchers have “proved” some hypothesis. But statistical analysis can never “prove” a hypothesis, in the sense of demonstrating

that it must be true (as one would in a logical or mathematical proof). Statistics can provide us with evidence, but it's always tentative and subject to the uncertainty that is ever present in the real world.

Sampling from a Population

The concept of aggregation implies that we can make useful insights by collapsing across data—but how much data do we need? The idea of *sampling* says that we can summarize an entire population based on just a small number of samples from the population, as long as those samples are obtained in the right way. For example, the PURE study enrolled a sample of about 135,000 people, but its goal was to provide insights about the billions of humans who make up the population from which those people were sampled. As we already discussed above, the way that the study sample is obtained is critical, as it determines how broadly we can generalize the results. Another fundamental insight about sampling is that, while larger samples are always better (in terms of their ability to accurately represent the entire population), there are diminishing returns as the sample gets larger. In fact, the rate at which the benefit of larger samples decreases follows a simple mathematical rule, growing as the square root of the sample size, such that in order to double the precision of our estimate we need to quadruple the size of our sample.

Causality and Statistics

The PURE study seemed to provide pretty strong evidence for a positive relationship between eating saturated fat and living longer, but this doesn't tell us what we really want to know: If we eat more saturated fat, will that cause us to live longer? This is because we don't know whether there is a direct causal relationship between eating saturated fat and living longer. The data are consistent with such a relationship, but they are equally consistent with some other factor causing both higher saturated fat and longer life. For example, one might imagine that people who are richer eat more saturated fat and richer people tend to live longer, but their longer life is not necessarily due to fat intake—it could instead be due to better health care, reduced psychological stress, better food quality, or many other factors. The PURE study investigators tried to account for these factors, but we can't be certain that their efforts completely removed the effects of other variables. The fact that other factors may explain the relationship between saturated fat intake and death is an example of why introductory statistics classes often teach that “correlation does not imply causation,” though the renowned data visualization expert Edward Tufte has added, “but it sure is a hint.”

Although observational research (like the PURE study) cannot conclusively demonstrate causal relationships, we generally think that causation can be demonstrated using studies that experimentally control and manipulate a specific factor. In medicine, such a study is referred to as a *randomized controlled trial* (RCT). Let's say that we wanted to do an

RCT to examine whether increasing saturated fat intake increases life span. To do this, we would sample a group of people and then assign them to either a treatment group (which would be told to increase their saturated fat intake) or a control group (who would be told to keep eating the same as before). It is essential that we assign the individuals to these groups randomly. Otherwise, people who choose the treatment group might be different in some way than people who choose the control group—for example, they might be more likely to engage in other healthy behaviors as well. We would then follow the participants over time and see how many people in each group died. Because we randomized the participants to treatment or control groups, we can be reasonably confident that there are no other differences between the groups that would *confound* the treatment effect; however, we still can't be certain because sometimes randomization yields treatment groups versus control groups that *do* vary in some important way. Researchers often try to address these confounds using statistical analyses, but removing the influence of a confound from the data can be very difficult.

A number of RCTs have examined the question of whether changing saturated fat intake results in better health and longer life. These trials have focused on *reducing* saturated fat because of the strong dogma among nutrition researchers that saturated fat is deadly; most of these researchers would have probably argued that it was not ethical to cause people to eat *more* saturated fat! However, the RCTs have shown a very consistent pattern: overall there is no appreciable effect on death rates of reducing saturated fat intake.

Suggested Reading

- *The Seven Pillars of Statistical Wisdom*, by Stephen Stigler. Stigler is one of the world's leading historians of statistics, and this book outlines what he sees as a number of the foundational ideas that underlie statistical thinking. His ideas strongly influenced the basic ideas that are presented in this chapter.
- *The Lady Tasting Tea: How Statistics Revolutionized Science in the Twentieth Century*, by David Salsburg. This book provides a readable yet detailed overview of the history of statistics, with a strong focus on amplifying the often overlooked contributions of women in the history of statistics.
- *Naked Statistics: Stripping the Dread from the Data*, by Charles Wheelan. A very fun tour of the main ideas of statistics.

Problems

1. Describe how statistics can be thought of as a set of tools for learning from data.
2. What does it mean to sample from a population, and why is this useful?
3. Describe the concept of a *randomized controlled trial* and outline the reason that we think that such an experiment can provide information about the causal effect of a treatment.
4. The three things that statistics can do for us are to _____, _____, and _____.

5. Early in the COVID-19 pandemic there was an observational study that reported effectiveness of the drug hydroxychloroquine in treating the disease. Subsequent randomized controlled trials showed no effectiveness of the drug for treating the disease. What might explain this discrepancy between observational results and randomized controlled trials? Choose all that apply.
- There were systematic differences in the observational study between those prescribed the drug and those who were prescribed other treatments.
 - Randomization helps eliminate differences between the treatment and control groups.
 - The drug is more effective when the physician gives it to the right patients.
6. Match the following examples with the most appropriate concept from the following list: aggregation, uncertainty, sampling.
- A researcher summarizes the scores of 10,000 people in a set of 12 numbers.
 - A person smokes heavily for 70 years but remains perfectly healthy and fit.
 - A researcher generalizes from a study of 1000 individuals to all humans.

INDEX

- alternative hypothesis, 94–96, 109, 110, 119, 230
analysis of variance, 170, 175, 178, 187, 189, 242, 245
- Bayes factor, 133–136, 143, 144, 157, 180, 183, 186, 230, 236, 247
Bayes' theorem, 61, 73–75, 78, 79
Bonferroni correction, 108–110
bootstrap, 86, 90–92, 115
- causality, 1, 5, 145, 152–154, 217
central limit theorem, 25, 84, 85, 213
chartjunk, 33, 39
clustered data, 214, 233, 239, 243, 247
clustering, 190, 200, 206, 210, 221
 hierarchical, 199
 K-means, 196–199, 210
coefficient of determination, 165, 166
Cohen's d , 116, 120, 122, 123
collider bias, 218, 219, 247
confidence interval, 90, 111–113, 115, 122–124, 130, 229, 235, 241
confounder, 216, 217
contingency table, 138, 140, 142–145
correlation coefficient, 123, 147–157, 159, 163, 174, 192, 193, 196, 201, 205, 207, 208
covariance, 148, 149, 154, 155, 163, 206, 207
credible interval, 124, 130, 137
cross-validation, 158, 173, 175, 258
- degrees of freedom, 52, 53, 97, 98, 100, 113, 139, 141, 145, 149, 164, 180, 187
dependent variable, 158, 160, 164, 171, 173–175, 190, 191, 212, 213, 216, 219, 222, 224, 231, 233, 238, 242
diagnostics, 171, 211, 212, 226, 227, 229, 234, 239, 243
dimensionality reduction, 190, 191, 200
- distribution
 binomial, 68, 88, 99, 125, 127, 131, 133, 179
 chi-squared, 139, 141, 145
 cumulative, 21, 27, 69, 88, 171
 frequency, 17, 27
 long-tailed, 17, 27, 213, 216, 242
 normal, 17, 25–27, 56, 83–85, 88, 90, 92, 97, 110, 113, 145, 162, 163, 171
 probability, 68, 69, 87, 96–98
 sampling, 80–85, 90, 91
 t , 98, 99, 109, 110, 113, 149, 187
 uniform, 88
dummy coding, 166, 167, 181, 182, 187, 229, 234
- effect size, 76, 107, 111, 116–118, 121–123, 182, 212, 224, 229, 235, 241, 254
Efron, Bradley, 90, 91
error
 sampling, 80, 81, 83, 91, 92
 Type I, 106, 110, 119, 120, 123, 227, 245
 Type II, 106, 110, 119, 145
Euclidean distance, 194–197, 201
eugenics, 161, 162
- factor analysis, 191, 204–206, 210
Fisher, Ronald, 104, 105, 161, 162, 241
- Galton, Francis, 160–162
- HARKing, 253, 257
histogram, 21–24, 26, 27, 43, 53, 82, 102, 137, 184, 192, 252
hypothesis testing, 93–95, 105, 107, 109, 124, 132, 138, 227, 239
- independence, 61, 66, 71, 78, 141, 145, 222, 238

- independent variable, 158–160, 175, 176, 190, 216, 217, 219, 224, 225, 227, 235, 239
Ioannidis, John, 250, 251, 257
- latent variable, 125, 152, 205–207
law of large numbers, 61, 63
leverage, 90, 214–216
lie factor, 35, 38
linear model, 47, 158, 165, 170, 171, 175–177, 182, 189, 222, 243
logit, 225, 235
- machine learning, 154, 173, 174
 unsupervised, 190, 191
measurement, 9, 11, 12, 14, 15, 46, 81, 83, 85, 116, 174, 183, 184, 204–206, 215
median, 3, 16, 30, 51, 59, 238
mixed-effects model, 245
mode, 16, 42, 44, 45, 52, 59
Monte Carlo simulation, 86, 87, 91, 92
multiple testing, 93, 107
multivariate statistics, 190–192, 200, 210
- NHANES, 18, 20–24, 30, 41, 44, 45, 50, 53, 70, 72, 82–84, 90, 95, 96, 111, 113, 115, 117, 143, 158, 172, 173, 178–180, 183
null hypothesis statistical testing, 93–95, 106, 109, 124, 125, 133, 136
- odds, 75, 77, 118, 119, 142, 226, 227, 229, 230, 235, 241
odds ratio, 76, 118, 119, 123, 142, 229, 230, 235, 236
outliers, 30, 32, 50, 51, 60, 150, 151, 157, 175, 214–216, 227, 247
overdispersion, 227, 234
overfitting, 48, 49, 59, 172, 175
- Pascal, Blaise, 64, 66
Pearson, Karl, 161, 162
p-hacking, 253–255, 257
positive predictive value, 248, 250
power analysis, 121–123
preregistration, 101, 230, 255
principal component analysis, 191, 201, 204, 210
probability, 61–79, 86, 95, 99, 101, 105, 106, 109, 112, 125, 127, 128, 133, 136, 138, 156, 250, 251
 classical, 64
 conditional, 61, 69–74, 78, 79, 125, 126
 empirical, 61, 63, 77
 pseudorandom numbers, 87, 92
 p-value, 93, 94, 99–106, 109, 110, 126, 140, 141, 150, 165, 180, 182, 188, 249, 256
 questionable research practices, 117, 252, 255
 randomization, 1, 7, 8, 99, 101–104, 109, 110
 randomized controlled trial, 6–8, 94, 106, 174
 regression
 linear, 158, 160, 167, 172, 174–176, 182, 190, 191, 201, 213, 216, 217, 224–227, 233, 238, 246, 247
 logistic, 213, 224–227, 229, 232, 247
 to the mean, 162, 163, 174
 reliability, 9, 12, 14
 replication, 106, 230, 250, 256
 representative sample, 4, 81, 85
 reproducibility, 248, 250, 254–258
- sample size, 6, 53, 63, 82, 85, 98, 100, 110, 112–116, 120–123, 148, 160, 208, 213, 239, 243
sampling, 6, 8, 49, 80, 81, 83, 97, 112, 179
 with replacement, 81, 91
 without replacement, 81, 91
scatterplot of matrices, 192, 210, 221
selection bias, 219
sign test, 178, 184, 185
Simpson's paradox, 138, 144, 145
standard deviation, 25, 52, 59, 82, 83, 90, 97, 112, 116, 122, 189, 246
standard error, 91, 97, 112, 114, 164
 of the mean, 80, 82, 83, 85, 90, 116, 122
 regression, 164
standardized residuals, 141, 142, 145
standardized scores, 56
statistical power, 111, 119–123, 175, 250–252, 256, 257
statistical significance, 94, 95, 104–108, 116, 250
- t* test, 98, 100, 101, 103, 110, 150, 175, 181, 182, 187, 223
 paired, 175, 185, 214
Tufte, Edward, 6, 28, 32, 34, 35, 37
- validity, 9, 12–14
variance, 52, 53, 97, 100, 122, 148, 154, 164, 165, 175, 187, 189, 201, 210, 246
winner's curse, 252, 257
Z-score, 40, 53, 55–57, 60