# STATS

## data & models

### FIFTH EDITION

De Veaux | Velleman | Bock
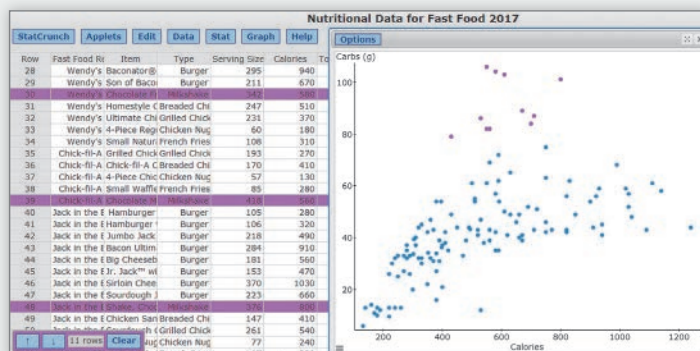
## Enrich assignments with question libraries

MyLab Statistics includes a number of question libraries providing additional opportunities for students to practice statistical thinking.

- StatCrunch Projects provide students with opportunities to analyze and interpret data. Each project consists of a series of questions about a large data set in StatCrunch.

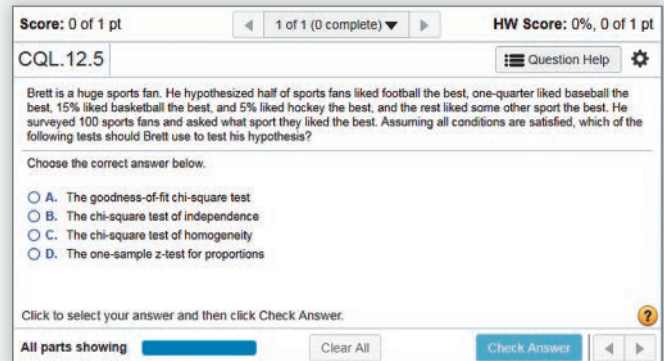- The Conceptual Question Library offers 1,000 conceptual-based questions to help students internalize concepts, make interpretations, and think critically about statistics.

| Score: 0 of 1 pt | ◀ 1 of 1 (0 complete) ▼ ▷ | HW Score: 0%, 0 of 1 pt |
|---|---|---|

CQL.12.5        ☰ Question Help ⚙

Brett is a huge sports fan. He hypothesized half of sports fans liked football the best, one-quarter liked baseball the best, 15% liked basketball the best, and 5% liked hockey the best, and the rest liked some other sport the best. He surveyed 100 sports fans and asked what sport they liked the best. Assuming all conditions are satisfied, which of the following tests should Brett use to test his hypothesis?

Choose the correct answer below.

- ○ A. The goodness-of-fit chi-square test
- ○ B. The chi-square test of independence
- ○ C. The chi-square test of homogeneity
- ○ D. The one-sample z-test for proportions

Click to select your answer and then click Check Answer.

All parts showing ▇▇▇▇▇   Clear All   Check Answer ◀ ▷

- The Getting Ready for Statistics Library contains more than 450 exercises on prerequisite topics. Assign these questions to students who may need a little extra practice on their prerequisite skills to be successful in your course.

- The StatTalk Video Library is based on a series of 24 videos, hosted by fun-loving statistician Andrew Vickers, that demonstrate important statistical concepts through interesting stories and real-life events.

## Incorporate additional author-created resources in class

Authors infuse their own voice, approach, and experiences teaching statistics into additional text-specific resources, such as interactive applets, technology manuals, workbooks, and more. Check out the Preface to learn more about what's available for this specific title.

**FIFTH EDITION**

# Stats: Data and Models

## Richard D. De Veaux
Williams College

## Paul F. Velleman
Cornell University (Emeritus)

## David E. Bock
Cornell University
Ithaca High School (Retired)

Pearson

*To Sylvia, who has helped me in more ways than she'll ever know,*
*and to Nicholas, Scyrine, Frederick, and Alexandra,*
*who make me so proud in everything that they are and do*
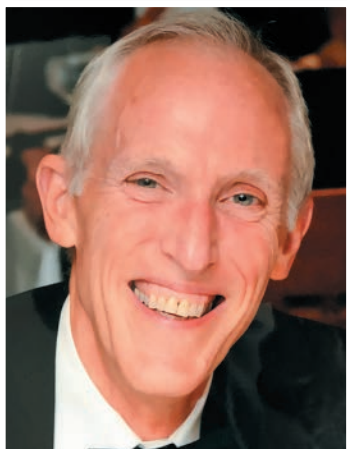
*—Dick*


*To my sons, David and Zev, from whom I've learned so much,*
*and to my wife, Sue, for taking a chance on me*

*—Paul*


*To Greg and Becca, great fun as kids and great friends as adults,*
*and especially to my wife and best friend, Joanna, for her*
*understanding, encouragement, and love*

*—Dave*

**Richard D. De Veaux** (Ph.D. Stanford University) is an internationally known educator, consultant, and lecturer. Dick has taught statistics at a business school (Wharton), an engineering school (Princeton), and a liberal arts college (Williams). While at Princeton, he won a Lifetime Award for Dedication and Excellence in Teaching. Since 1994, he has taught at Williams College, although he returned to Princeton for the academic year 2006–2007 as the William R. Kenan Jr. Visiting Professor of Distinguished Teaching. He is currently the C. Carlisle and Margaret Tippit Professor of Statistics at Williams College. Dick holds degrees from Princeton University in Civil Engineering and Mathematics and from Stanford University, where he studied statistics with Persi Diaconis and dance with Inga Weiss. His research focuses on the analysis of large datasets and data mining in science and industry. Dick has won both the Wilcoxon and Shewell awards from the American Society for Quality. He is an elected member of the International Statistics Institute (ISI) and a Fellow of the American Statistical Association (ASA). Dick was elected Vice President of the ASA in 2018 and will serve from 2019 to 2021. Dick is also well known in industry, having consulted for such *Fortune* 500 companies as American Express, Hewlett-Packard, Alcoa, DuPont, Pillsbury, General Electric, and Chemical Bank. He was named the "Statistician of the Year" for 2008 by the Boston Chapter of the American Statistical Association. In his spare time he is an avid cyclist and swimmer, and is a frequent singer and soloist with various local choirs, including the Choeur Vittoria of Paris, France. Dick is the father of four children.

**Paul F. Velleman** (Ph.D. Princeton University) has an international reputation for innovative statistics education. He designed the Data Desk® software package and is also the author and designer of the award-winning ActivStats® multimedia software, for which he received the EDUCOM Medal for innovative uses of computers in teaching statistics and the ICTCM Award for Innovation in Using Technology in College Mathematics. He is the founder and CEO of Data Description, Inc. (**www.datadesk.com**), which supports both of these programs. Data Description also developed and maintains the Internet site *Data and Story Library* (DASL; **dasl.datadescription.com**), which provides all of the datasets used in this text as well as many others useful for teaching statistics, and the statistics conceptual tools at **astools.datadesk.com**. Paul coauthored (with David Hoaglin) the book *ABCs of Exploratory Data Analysis*. Paul is Emeritus Professor of Statistical Sciences at Cornell University, where he was awarded the MacIntyre Prize for Exemplary Teaching. Paul earned his M.S. and Ph.D. from Princeton University, where he studied with John Tukey. His research often focuses on statistical graphics and data analysis methods. Paul is a Fellow of the American Statistical Association and of the American Association for the Advancement of Science. He was a member of the working group that developed the GAISE 2016 guidelines for teaching statistics.

**David E. Bock** taught mathematics at Ithaca High School for 35 years. He has taught Statistics at Ithaca High School, Tompkins-Cortland Community College, Ithaca College, and Cornell University. Dave has won numerous teaching awards, including the MAA's Edyth May Sliffe Award for Distinguished High School Mathematics Teaching (twice), Cornell University's Outstanding Educator Award (three times), and has been a finalist for New York State Teacher of the Year.

Dave holds degrees from the University at Albany in Mathematics (B.A.) and Statistics/Education (M.S.). Dave has been a reader and table leader for the AP Statistics exam and a Statistics consultant to the College Board, leading workshops and institutes for AP Statistics teachers. His understanding of how students learn informs much of this book's approach.

Richard De Veaux, Paul Velleman, and David Bock have authored several successful books in the introductory college and AP High School market including *Intro Stats*, Fifth Edition (Pearson, 2018) and *Stats: Modeling the World*, Fifth Edition (Pearson, 2019).

# TABLE OF CONTENTS

---

*Indicates optional sections.

## PART III  Gathering Data

## PART IV  Randomness and Probability

*S*tats: Data and Models, fifth edition, has been especially exciting to develop. The book you hold steps beyond our previous editions in several important ways. Of course, we've kept our conversational style and anecdotes,[1] but we've enriched that material with tools for teaching about randomness, sampling distribution models, and inference throughout the book. And we've expanded discussions of models for data to introduce models with more than two variables earlier in the text. We've taken our inspiration both from our experience in the classroom and from the 2016 revision of the Guidelines for Assessment and Instruction in Statistics Education (GAISE) report adopted by the American Statistical Association. As a result, we increased the text's innovative uses of technology to encourage more statistical thinking, while maintaining its traditional core concepts and coverage. You'll notice that, to expand our attention beyond just one or two variables, we've adjusted the order of some topics.

# Innovations

## Technology

One of the new GAISE guidelines states: *Use technology to explore concepts and analyze data.* We think a modern statistics text should recognize from the start that statistics is practiced with technology. And so should our students. You won't find tedious calculations worked by hand. You *will* find equation forms that favor intuition over calculation. You'll find extensive use of real data—even large datasets. Throughout, you'll find a focus on statistical thinking rather than calculation. The question that motivates each of our hundreds of examples is not "How do you calculate the answer?" but "How do you think about the answer?"

For this edition of *Stats: Data and Models* we've taken this principle still further. We have harnessed technology to improve the learning of two of the most difficult concepts in the introductory course: the idea of a sampling distribution and the reasoning of statistical inference.

## Multivariable Thinking and Multiple Regression

GAISE's first guideline is to give students experience with multivariable thinking. The world is not univariate, and relationships are not limited to two variables. This edition of *Stats: Data and Models* introduces a third variable as early as Chapter 3's discussion of contingency tables and mosaic plots. Then, following the discussion of correlation and regression as a tool (that is, without inference) in Chapters 6, 7, and 8, we introduce multiple regression in Chapter 9.

Multiple regression may be the most widely used statistical method, and it is certainly one that students need to understand. It is easy to perform multiple regressions with any statistics program, and the exercise of thinking about more than two variables early in the course is worth the effort. We've added new material about interpreting what regression models say. The effectiveness of multiple regression is immediately obvious and makes the reach and power of statistics clear. The use of real data underscores the universal applicability of these methods.

When we return to regression in Chapters 23 and 24 to discuss inference, we can deal with both simple and multiple regression models together. There is nothing different to discuss. (For this reason we set aside the *F*-test until the chapter on ANOVA.)

---

[1]And footnotes

Innovative ways to teach the logic of statistical inference have received increasing attention. Among these are greater use of computer-based simulations and resampling methods (randomization tests and bootstrapping) to teach concepts of inference.

## Bootstrap

The introduction to the new GAISE guidelines explicitly mentions the bootstrap method. The bootstrap is not as widely available or as widely understood as multiple regression. But it follows our presentation naturally. In this edition, we introduce a new feature, **Random Matters**. Random Matters elements in early chapters draw small samples repeatedly from large populations to illustrate how the randomness introduced by sampling leads to both sampling distributions and statistical reasoning for inference. But what can we do when we have only a sample? The bootstrap provides a way to continue this line of thought, now by resampling from the sample at hand.

Bootstrapping provides an elegant way to simulate sampling distributions that we might not otherwise be able to see. And it does not require the assumption of Normality expected by Student's *t*-based methods. However, these methods are not as widely available or widely used in other disciplines, so they should not be the only—or even the principal—methods taught. They may be able to enhance student understanding, but instructors may wish to downplay them if that seems best for a class. We've placed these sections strategically so that instructors can choose the level that they are comfortable with and that works best with their course.

## Real Data

GAISE recommends that instructors integrate real data with a context and purpose. More and more high school math teachers are using examples from statistics to demonstrate intuitively how a little bit of math can help us say a lot about the world. So our readers expect statistics to be about real-world insights. *Intro Stats* keeps readers engaged and interested because we show statistics in action right from the start. The exercises pose problems of the kind likely to be encountered in real life and propose ways to think about making inferences almost immediately—and, of course, always with real, up-to-date data.

Let us be clear. *Stats: Data and Models* comes with an archive of more than 500 datasets used in more than 700 applications throughout the book. The datasets are available online at the student resource site, in the DASL archive and in MyLab Statistics. Examples that use these datasets cite them in the text. Exercises are marked when they use one of them; exercise names usually indicate the name of the dataset. We encourage students to get the datasets and reproduce our examples using their statistics software, and some of the exercises require that.

## Streamlined Content

Following the GAISE recommendations, we've streamlined several parts of the course: Introductory material is covered more rapidly. Today's students have seen a lot of statistics in their K–12 math courses and in their daily contact with online and news sources. We still cover the topics to establish consistent terminology (such as the difference between a histogram and a bar chart). Chapter 2 does most of the work that previously took two chapters.

The Random Matters features show students that statistics vary from sample to sample, show them (empirical) sampling distributions, note the effect of sample size on the shape and variation of the sampling distribution of the mean, and suggest that it looks Normal. As a result, the discussion of the Central Limit Theorem is transformed from the most difficult one in the course to a relatively short discussion ("What you think is true about means really is true; there's this theorem.") that can lead directly to the reasoning of confidence intervals.

Finally, introducing multiple regression doesn't really add much to the lesson on inference for multiple regression because little is new.

## GAISE 2016

As we've said, all of these enhancements follow the new Guidelines for Assessment and Instruction in Statistics Education (GAISE) 2016 report adopted by the American Statistical Association:

1. Teach statistical thinking.
   ◆ Teach statistics as an investigative process of problem solving and decision making.
   ◆ Give students experience with multivariable thinking.
2. Focus on conceptual understanding.
3. Integrate real data with a context and purpose.
4. Foster active learning.
5. Use technology to explore concepts and analyze data.
6. Use assessments to improve and evaluate student learning.

The result is a course that is more aligned with the skills needed in the 21st century, one that focuses even more on statistical thinking and makes use of technology in innovative ways, while retaining core principles and topic coverage.

The challenge has been to use this modern point of view to improve learning without discarding what is valuable in the traditional introductory course. Many first statistics courses serve wide audiences of students who need these skills for their own work in disciplines where traditional statistical methods are, well, traditional. So we have not reduced our emphasis on the concepts and methods you expect to find in our texts.

# Chapter Order

We've streamlined the presentation of basic topics that most students have already seen. Pie charts, bar charts, histograms, and summary statistics all appear in Chapter 2. Chapter 3 introduces contingency tables, and Chapter 4 discusses comparing distributions. Chapter 5 introduces the Normal model and the 68–95–99.7 Rule. The four chapters of Part II then explore linear relationships among quantitative variables—but here we introduce only the models and how they help us understand relationships. We leave the inference questions until later in the book. Part III discusses how data are gathered by survey and experiment.

Part IV provides background material on probability, random variables, and probability models. In Part V, Chapter 16 introduces confidence intervals for proportions as soon as we've reassured students that their intuition about the sampling distribution of proportions is correct. Chapter 17 formalizes the Central Limit Theorem and introduces Student's $t$-models. Chapter 18 is then about testing hypotheses, and Chapter 19 elaborates further, discussing alpha levels, Type I and Type II errors, power, and effect size. The chapters in Part VI deal with comparing groups (both with proportions and with means), paired samples, chi-square. Finally, Part VII discusses inferences for regression models (both simple and multiple), intelligent uses of multiple regression, and Analysis of Variance, both one- and two-way. A final chapter on data mining looks to the future.

We've found that one of the challenges students face is how to know what technique to use when. In the real world, questions don't come at the ends of the chapters. So, as always, we've provided summaries at the end of each part along with a series of exercises designed to stretch student understanding. These Part Reviews are a mix of questions from all the chapters in that part. The final set are "book-level" review problems that ask students to integrate what they've learned from the entire course. The questions range

from simple questions about what method to use in various situations to a more complete data analyses from real data. We hope that these will provide a useful way for students to organize their understanding at the end of the course.

# Our Approach

We've discussed how this book is different, but there are some things we haven't changed.

◆ *Readability.* This book doesn't read like other statistics texts. Our style is both colloquial and informative, engaging students to actually read the book to see what it says.

◆ *Humor.* You will find quips and wry comments throughout the narrative, in margin notes, and in footnotes.

◆ *Informality.* Our informal diction doesn't mean that we treat the subject matter lightly or informally. We try to be precise and, wherever possible, we offer deeper explanations and justifications than those found in most introductory texts.

◆ *Focused lessons.* The chapters are shorter than in most other texts so that instructors and students can focus on one topic at a time.

◆ *Consistency.* We try to avoid the "do what we say, not what we do" trap. Having taught the importance of plotting data and checking assumptions and conditions, we model that behavior through the rest of the book. (Check out the exercises in Chapter 24.)

◆ *The need to read.* Statistics is a consistent story about how to understand the world when we have data. The story can't be told piecemeal. This is a book that needs to be read, so we've tried to make the reading experience enjoyable. Students who start with the exercises and then search back for a worked example that looks the same but with different numbers will find that our presentation doesn't support that approach.

## Mathematics

Mathematics can make discussions of statistics concepts, probability, and inference clear and concise. We don't shy away from using math where it can clarify without intimidating. But we know that some students are discouraged by equations, so we always provide a verbal description and a numerical example as well.

Nor do we slide in the opposite direction and concentrate on calculation. Although statistics calculations are generally straightforward, they are also usually tedious. And, more to the point, today, virtually all statistics are calculated with technology. We have selected the equations that focus on illuminating concepts and methods rather than for hand calculation. We sometimes give an alternative formula, better suited for hand calculation, for those who find that following the calculation process is a better way to learn about the result.

## Technology and Data

We assume that computers and appropriate software are available—at least for demonstration purposes. We hope that students have access to computers and statistics software for their analyses.

We discuss generic computer output at the end of most chapters, but we don't adopt any particular statistics software. The **Tech Support** sections at the ends of chapters offer guidance for seven common software platforms: Data Desk, Excel, JMP, Minitab, SPSS, StatCrunch, and R. We also offer some advice for TI-83/84 Plus graphing calculators, although we hope that those who use them will also have some access to computers and statistics software.

We don't limit ourselves to small, artificial datasets, but base most examples and exercises on real data with a moderate number of cases. Machine-readable versions of the data are available at the book's website, pearsonhighered.com/dvb and at dasl.datadescription.com.

# Features

## Enhancing Understanding

**Where Are We Going?** Each chapter starts with a paragraph that raises the kinds of questions we deal with in the chapter. A chapter outline organizes the major topics and sections.

**New! Random Matters.** This new feature travels along a progressive path of understanding randomness and our data. The first Random Matters element begins our thinking about drawing inferences from data. Subsequent Random Matters draw histograms of sample means, introduce the thinking involved in permutation tests, and encourage judgment about how likely the observed statistic seems when viewed against the simulated sampling distribution of the null hypothesis (without, of course, using those terms).

**Margin and in-text boxed notes.** Throughout each chapter, boxed margin and in-text notes enhance and enrich the text.

**Reality Check.** We regularly remind students that statistics is about understanding the world with data. Results that make no sense are probably wrong, no matter how carefully we think we did the calculations. Mistakes are often easy to spot with a little thought, so we ask students to stop for a reality check before interpreting their result.

**Notation Alert.** Throughout this book, we emphasize the importance of clear communication, and proper notation is part of the vocabulary of statistics. We've found that it helps students when we are clear about the letters and symbols statisticians use to mean very specific things, so we've included Notation Alerts whenever we introduce a special notation that students will see again.

Each chapter ends with several elements to help students study and consolidate what they've seen in the chapter.

◆ **What Can Go Wrong?** sections highlight the most common errors that people make and the misconceptions they have about statistics. One of our goals is to arm students with the tools to detect statistical errors and to offer practice in debunking misuses of statistics, whether intentional or not.

◆ **Connections** specifically ties the new topics to those learned in previous chapters.

◆ The **Chapter Review** summarizes the story told by the chapter and provides a bullet list of the major concepts and principles covered.

◆ A **Review of Terms** is a glossary of all of the special terms introduced in the chapter. In the text, these are printed in **bold** and underlined. The Review provides page references, so students can easily turn back to a full discussion of the term if the brief definition isn't sufficient.

The **Tech Support** section provides the commands in each of the supported statistics packages that deal with the topic covered by the chapter. These are not full documentation, but should be enough to get a student started in the right direction.

## Learning by Example

**Step-by-Step Examples.** We have expanded and updated the examples in our innovative Step-by-Step feature. Each one provides a longer, worked example that guides students through the process of analyzing a problem. The examples follow our three-step Think, Show, Tell organization for approaching a statistics task. They are organized with general explanations of each step on the left and a worked-out solution on the right. The right side of the grid models what would be an "A" level solution to the problem. Step-by-Steps illustrate the importance of thinking about a statistics question (What do we know? What do we hope to learn? Are the assumptions and conditions satisfied?) and reporting our findings (the Tell step). The Show step contains the mechanics of calculating results and conveys our belief that it is only one part of the process. Our emphasis is on statistical thinking, and the pedagogical result is a better understanding of the concept, not just number crunching.

**Examples.** As we introduce each important concept, we provide a focused example that applies it—usually with real, up-to-the-minute data. Many examples carry the discussion through the chapter, picking up the story and moving it forward as students learn more about the topic.

**Just Checking.** Just Checking questions are quick checks throughout the chapter; most involve very little calculation. These questions encourage students to pause and think about what they've just read. The Just Checking answers are at the end of the exercise sets in each chapter so students can easily check themselves.

## Assessing Understanding

Our **Exercises** have some special features worth noting. First, you'll find relatively simple, focused exercises organized by chapter section. After that come more extensive exercises that may deal with topics from several parts of the chapter or even from previous chapters as they combine with the topics of the chapter at hand. All exercises appear in pairs. The odd-numbered exercises have answers in the back of student texts. Each even-numbered exercise hits the same topic (although not in exactly the same way) as the previous odd exercise. But the even-numbered answers are not provided. If a student is stuck on an even exercise, looking at the previous odd one (and its answer) can often provide the help needed.

More than 600 of our exercises have a 🅣 tag next to them to indicate that the dataset referenced in the exercise is available electronically. The exercise title or a note provides the dataset title. Some exercises have a 🎲 tag to indicate that they call for the student to generate random samples or use randomization methods such as the bootstrap. Although we hope students will have access to computers, we provide ample exercises with full computer output for students to read, interpret, and explain.

We place all the exercises—including section-level exercises—at the end of the chapter. Our writing style is colloquial and encourages reading. We are telling a story about how to understand the world when you have data. Interrupting that story with exercises every few pages would encourage a focus on the calculations rather than the concepts.

**Part Reviews.** The book is partitioned into seven conceptual parts; each ends with a Part Review. The part review discusses the concepts in that part of the text, tying them together and summarizing the story thus far. Then there are more exercises. These exercises have the advantage (for study purposes) of not being tied to a chapter, so they lack the hints of what to do that would come from that identification. That makes them more like potential exam questions and a good tool for review. Unlike the chapter exercises, these are not paired.

**Parts I-VII Cumulative Review Exercises.** Cumulative Review exercises are longer and cover concepts from the book as a whole.

# Additional Resources Online

Most of the supporting materials can be found online:

At the book's website at **pearsonhighered.com/dvb**

Within the MyLab Statistics course at **pearson.com/mylab/statistics**

Datasets are also available at **dasl.datadesk.com**.

**Data desk** is a statistics program with a graphical interface that is easy to learn and use. A student version is available at **datadesk.com**. Click on the **Teachers & Students** tab at the top of the page.

New tools that provide interactive versions of the distribution tables at the back of the book and tools for randomization inference methods such as the bootstrap and for repeated sampling from larger populations can be found online at **astools.datadesk.com**.

# MyLab Statistics for *Stats: Data & Models, 5e*

(access code required)

MyLab Statistics is available to accompany Pearson's market-leading text offerings. To give students a consistent tone, voice, and teaching method, each text's flavor and approach are tightly integrated throughout the accompanying MyLab course, making learning the material as seamless as possible.

### NEW! StatCrunch Projects

StatCrunch Projects provide opportunities for students to explore data beyond the classroom. In each project, students analyze a data set in StatCrunch® and answer assignable MyLab questions for immediate feedback. StatCrunch Projects span the entire curriculum or focus on certain key concepts. Questions from each project can also be assigned individually.



### UPDATED! Real-World Data

Statistical concepts are applied to everyday life through the extensive current, real-world data examples and exercises provided throughout the text.



### NEW! Interactive Applets

Author-created interactive applets take an experiential approach to helping students learn statistical concepts. They are available in MyLab Statistics and at astools.datadesk.com.

# Resources for Success

## Instructor Resources

### Annotated Instructor's Edition
Includes answers to all text exercises, as well as a set of Instructor Resource Pages that offer chapter-by-chapter teaching suggestions and commentary. (ISBN-13: 978-0-13-516396-2; ISBN-10: 0-13-516396-X)

### Instructor's Solutions Manual (Download Only)
This manual contains detailed solutions to all of the exercises. These files can be downloaded from within MyLab Statistics or from **www.pearson.com**.

### Instructor Resource Guide (Download Only)
This resource guide includes chapter-by-chapter comments on the major concepts, tips on presenting topics, extra teaching examples, and sample quizzes, tests, and projects. These files can be downloaded from within MyLab Statistics or from **www.pearson.com**.

### TestGen
TestGen® (www.pearson.com/testgen) enables instructors to build, edit, print, and administer tests using a computerized bank of questions developed to cover all the objectives of the text. TestGen is algorithmically based, allowing instructors to create multiple but equivalent versions of the same question or test with the click of a button. Instructors can also modify test bank questions or add new questions. The software and test bank are available for download from Pearson's online catalog, **www.pearson.com**. The questions are also assignable in MyLab Statistics.

### PowerPoint Lecture Slides
PowerPoint Lecture Slides provide an overview of each chapter, stressing important definitions and offering additional examples. They can be downloaded from MyLab Statistics or from **www.pearson.com**.

### Learning Catalytics
Now included in all MyLab Statistics courses, this student response tool uses students' smartphones, tablets, or laptops to engage them in more interactive tasks and thinking during lecture. Learning Catalytics™ fosters student engagement and peer-to-peer learning with real-time analytics. Access pre-built exercises created specifically for statistics.

### Question Libraries
In addition to Statcrunch Projects, MyLab Statistics also includes a Getting Ready for Statistics library that contains more than 450 exercises on prerequisite topics and a Conceptual
Question Library with 1000 questions that assess conceptual understanding.

### Minitab and Minitab Express™
Bundling Minitab and Minitab Express software with the text ensures students have access to to the software they need for the duration of their course. ISBN 13: 978-0-13-445640-9; ISBN 10: 0-13-445640-8

### JMP Student Edition
An easy-to-use, streamlined version of JMP desktop statistical discovery software from SAS Institute, Inc. is available for bundling with the text. ISBN-13: 978-0-13-467979-2; ISBN-10: 0-13-467979-2

### XLSTAT™
An Excel add-in that enhances the analytical capabilities of Excel, XLSTAT is used by leading businesses and universities around the world. It is available to bundle with this text. For more information go to **www.pearsonhighered.com/xlstatupdate**. ISBN-13: 978-0-321-75932-0; ISBN-10: 0-321-75932-X

### Accessibility
Pearson works continuously to ensure our products are as accessible as possible to all students. We are working toward achieving WCAG 2.0 Level AA and Section 508 standards, as expressed in the Pearson Guidelines for Accessible Educational Web Media, **www.pearson.com/mylab/statistics/accessibility**.

## Student Resources

### Video Resources
Step-by-Step Example videos guide students through the process of analyzing a problem using the "Think, Show, and Tell" strategy from the textbook. StatTalk Videos, hosted by fun-loving statistician Andrew Vickers, demonstrates important statistical concepts through interesting stories and real-life events. StatTalk videos come with accompanying MyLab assessment questions.

### StatCrunch
StatCrunch® is powerful web-based statistical software that allows users to collect, crunch, and communicate with data. The vibrant online community offers tens of thousands of shared datasets for students and instructors to analyze, in addition to all of the datasets in the text or online homework. StatCrunch is integrated directly into MyLab Statistics or it can be purchased separately. Learn more at **www.statcrunch.com**.

### Datasets Available Online
Data sets can be found at **pearsonhighered.com/dvb** and (DASL) at **www.DASL.datadescription.com** which holds all of the datasets for the book as well as many others. Search by dataset name to find datasets for the exercises in the book. Search by statistics method to find examples for lessons or for assignments. Datasets can be easily transferred to any statistics program.

### Statistical Software Support
Instructors and students can copy datasets from the text, DASL, and MyLab exercises directly into software such as StatCrunch, Data desk, or Excel®. Students can also access instructional support tools including tutorial videos, Study Cards, and manuals for a variety of statistical software programs including StatCrunch, Excel, Minitab®, JMP®, R, SPSS, and TI 83/84 calculators.

**pearson.com/mylab/statistics**

## Accounting

Double-checking procedures, 347

## Advertising

Product endorsements, 601
Radio ads, 594
Sex and violence, 376, 645, 667, 672
Super Bowl commercials, 368
TV ads, 569

## Agriculture

Apples, 466
Beetles, 371
Egg production, 599, 666–667
Farmers' markets, 441
Fertilizers, 829
Global climate change, 14
Milk production, 673
Oranges, 281, 437
Peas, 865–866
Seeds, 439, 568
Tomatoes, 156, 367, 369
Tree growth, 281
Vineyards, 14, 121, 314, 377

## Banking

Credit cards, 160, 362, 465, 494, 503, 530, 698
Customers' ages, 698
Loans, 593
Online banking, 414
Tellers, 467, 826

## Business (General)

Assets of corporations, 123, 124
Contracts, 438–439
Mergers, 471
Women executives, 568
Women-owned businesses, 472, 599, 601

## Company Names

A.G. Edwards, 97
Allstate Insurance Company, 534
Amazon, 2–3, 5, 871
Arby's, 14
Bentley, 253
Burger King, 200–201, 202–203, 205, 210, 211, 213–214, 216, 233, 269, 298, 762–764, 765–766, 788
Casualty Actuarial Society, 419
Cleveland Casting Plant, 14
*Consumer Reports,* 4, 7, 15, 640
Cornell University, 534
Daimler AG, 253
Dartmouth College, 529
Food and Drug Administration, 352
Guinness Company, 511
Hostess Company, 328
Lay's, 534
Mars, 388, 395, 699
McDonald's, 788
Nabisco Company, 570
National Strategy for Trusted Identities in Cyberspace, 6
Nissan, 15, 38–39
OkCupid, 64, 65, 78
Paralyzed Veterans of America, 503, 568, 867, 873, 874–875, 877–878, 882–883
Rolls-Royce, 253
Scripps Institution of Oceanography, 240–241
Sleep Foundation, 550–551
SmartWool, 539, 540
Society of Actuaries, 419
*Sports Illustrated,* 229
Texaco, Inc., 828
Toyota, 233
Verizon, 76–78
White Star Line, 34

## Consumers

Consumer attitudes, 345, 415
Credit card spending, 30, 35, 267
Grocery shopping, 13
Laundry detergent, 371, 375
Online shopping, 13, 84, 231
Wardrobe, 392

## Demographics

Adoptions, 54
Age and political party, 377
Age of athletes, 637
Age of bank customers, 698
Age of spouses, 656, 659
Birthrates, 234–235, 906
Consumer survey, 415
Deaths, 119–120
Fertility rates, 276
Foreign-born citizens, 636
Hispanics, 566
Life expectancy, 54
Marriage trends, 52, 120, 244–245, 269, 273, 274, 275, 500, 566, 599, 740, 741
Population growth, 62, 119
Poverty and region, 87
Religion, 396, 600
State populations, 60
Stay-at-home dads, 640

## Distribution and Operations Management

Delivery service, 93
Packaging stereos, 432–434
Refurbished computers, 423–424
Shipments, 151, 502

## E-Commerce

Earnings, 367
Online insurance, 668, 669
Online shopping, 13, 84, 231, 467, 468, 502
Website sales, 464, 539

## Economics

Boomtowns, 62
Cost of living, 29, 117–118, 125, 234
Earnings of college graduates, 736, 737
Earnings predictions, 738
GDP, 194–195, 278, 279
Global comparisons, 169, 175
Human Development Index, 270
Incomes, 151, 195, 230, 788–789
Inflation, 276
Interest rates, 195–196, 230, 273, 274, 278
Labor force participation, 86
Living conditions, 277, 278
Market segments, 267
Nest Egg Index, 97
Stock market, 315
Wealth redistribution, 756

## Education

Absenteeism, 568
ACT scores, 132, 151, 152
Age and educational attainment, 704
Birth order and college, 415, 416, 417
Cartoons and test performance, 87
Cheating on tests, 599
College admission rates, 80, 393–394
College attendance, 344
College graduation, 417
College homecoming, 376
College majors, 416
College meal plans, 472, 529, 601
College professors, 154
College retention rate, 504, 505
College tuition, 671
Computer lab fees, 532, 569–570
Computer software, 374
Cost, 748
Dorm amenities, 414

## Energy

## Environment

## Famous People

## Finance and Investments

## Pharmaceuticals, Medicine, and Health

## Politics and Popular Culture

# Stats Starts Here[1]

## WHERE ARE WE GOING?

Statistics gets no respect. People say things like "You can prove anything with statistics." People will write off a claim based on data as "just a statistical trick." And statistics courses don't have the reputation of being students' first choice for a fun elective.

But statistics *is* fun. That's probably not what you heard on the street, but it's true. Statistics is the science of learning from data. A little practice thinking statistically is all it takes to start seeing the world more clearly and accurately.

This is a text about understanding the world by using data. So we'd better start by understanding data. There's more to that than you might have thought.

> But where shall I begin?" asked Alice. "Begin at the beginning," the King said gravely, "and go on till you come to the end: then stop.
>
> —*Lewis Carroll,*
> Alice's Adventures
> in Wonderland

## 1.1 What Is Statistics?

People around the world have one thing in common—they all want to figure out what's going on. You'd think with the amount of information available to everyone today this would be an easy task, but actually, as the amount of information grows, so does our need to understand what it can tell us.

At the base of all this information, on the Internet and all around us, are data. We'll talk about data in more detail in the next section, but for now, think of **data** as any collection of numbers, characters, images, or other items that provide information about something. What sense can we make of all this data? You certainly can't make a coherent picture from random pieces of information. Whenever there are data and a need for understanding the world, you'll find statistics.

This text will help you develop the skills you need to understand and communicate the knowledge that can be learned from data. By thinking clearly about the question you're trying to answer and learning the statistical tools to show what the data are saying, you'll acquire the skills to tell clearly what it all means. Our job is to help you make sense of the concepts and methods of statistics and to turn it into a powerful, effective approach to understanding the world through data.

---

[1]We were thinking of calling this chapter "Introduction" but nobody reads the introduction, and we wanted you to read this. We feel safe admitting this down here in the footnotes because nobody reads footnotes either.

> **"** Data is king at Amazon. Clickstream and purchase data are the crown jewels at Amazon. They help us build features to personalize the Web site experience. **"**
>
> —*Ronny Kohavi,*
> *former Director of Data*
> *Mining and Personalization,*
> *Amazon.com*

*Q:* What is statistics?

*A:* Statistics is a way of reasoning, along with a collection of tools and methods, designed to help us understand the world.

*Q:* What are statistics?

*A:* Statistics (plural) are particular calculations made from data.

*Q:* So what is data?

*A:* You mean "what *are* data?" Data is the plural form. The singular is datum.

*Q:* OK, OK, so what are data?

*A:* Data are values along with their context.

The ads say, "Don't drink and drive; you don't want to be a statistic." But you can't be a statistic.

We say, "Don't be a datum."

Data vary. Ask different people the same question and you'll get a variety of answers. Statistics helps us to make sense of the world described by our data by seeing past the underlying variation to find patterns and relationships. This text will teach you skills to help with this task and ways of thinking about variation that are the foundation of sound reasoning about data.

Consider the following:

◆ If you have a Facebook account, you have probably noticed that the ads you see online tend to match your interests and activities. Coincidence? Hardly. According to *The Wall Street Journal* (10/18/2010),[2] much of your personal information has probably been sold to marketing or tracking companies. Why would Facebook give you a free account and let you upload as much as you want to its site? Because your data are valuable! Using your Facebook profile, a company might build a profile of your interests and activities: what movies and sports you like; your age, sex, education level, and hobbies; where you live; and, of course, who your friends are and what *they* like. From Facebook's point of view, your data are a potential gold mine. Gold ore in the ground is neither very useful nor pretty. But with skill, it can be turned into something both beautiful and valuable. What we're going to talk about is how you can mine your own data and learn valuable insights about the world.

◆ Americans spend an average of 4.9 hours per day on their smartphones. Trillions of text messages are sent each year.[3] Some of these messages are sent or read while the sender or the receiver is driving. How dangerous is texting while driving?

How can we study the effect of texting while driving? One way is to measure reaction times of drivers faced with an unexpected event while driving and texting. Researchers at the University of Utah tested drivers on simulators that could present emergency situations. They compared reaction times of sober drivers, drunk drivers, and texting drivers.[4] The results were striking. The texting drivers actually responded more slowly and were more dangerous than drivers who were above the legal limit for alcohol.

In this text, you'll learn how to design and analyze experiments like this. You'll learn how to interpret data and to communicate the message you see to others. You'll also learn how to spot deficiencies and weaknesses in conclusions drawn by others that you see in newspapers and on the Internet every day. Statistics can help you become a more informed citizen by giving you the tools to understand, question, and interpret data.

## 1.2 Data

*STATISTICS IS ABOUT . . .*

• Variation: Data vary because we don't see everything, and even what we do see, we measure imperfectly.

• Learning from data: We hope to learn about the world as best we can from the limited, imperfect data we have.

• Making intelligent decisions: The better we understand the world, the wiser our decisions will be.

Amazon.com opened for business in July 1995, billing itself as "Earth's Biggest Bookstore." By 1997, Amazon had a catalog of more than 2.5 million book titles and had sold books to more than 1.5 million customers in 150 countries. In 2017, the company's sales reached almost $178 billion (more than 30% over the previous year). Amazon has sold a wide variety of merchandise, including a $400,000 necklace, yak cheese from Tibet, and the largest book in the world. How did Amazon become so successful and how can it keep track of so many customers and such a wide variety of products? The answer to both questions is *data*.

But what are data? Think about it for a minute. What exactly *do* we mean by "data"? You might think that data have to be numbers, but data can be text, pictures, web pages,

---

[2]blogs.wsj.com/digits/2010/10/18/referers-how-facebook-apps-leak-user-ids/

[3]informatemi.com/blog/?p=133

[4]"Text Messaging During Simulated Driving," Drews, F. A., et al., Human Factors: hfs.sagepub.com/content/51/5/762

and even audio and video. If you can sense it, you can measure it. Data are now being collected automatically at such a rate that IBM estimates that "90% of the data in the world today has been created in the last two years alone."[5]

Let's look at some hypothetical values that Amazon might collect:

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| B0000010AA | 0.99 | Chris G. | 902 | 105-2686834-3759466 | 1.99 | 0.99 | Illinois |
| Los Angeles | Samuel R. | Ohio | N | B000068ZVQ | Amsterdam | New York, New York | Katherine H. |
| Katherine H. | 002-1663369-6638649 | Beverly Hills | N | N | 103-2628345-9238664 | 0.99 | Massachusetts |
| 312 | Monique D. | 105-9318443-4200264 | 413 | B00000I5Y6 | 440 | B000002BK9 | 0.99 |
| Canada | Detroit | 440 | 105-1372500-0198646 | N | B002MXA7Q0 | Ohio | Y |

Try to guess what they represent. Why is that hard? Because there is no *context*. If we don't know what values are measured and what is measured about them, the values are meaningless. We can make the meaning clear if we organize the values into a **data table** such as this one:

| Order Number | Name | State/Country | Price | Area Code | Download | Gift? | ASIN | Artist |
|---|---|---|---|---|---|---|---|---|
| 105-2686834-3759466 | Katherine H. | Ohio | 0.99 | 440 | Amsterdam | N | B0000015Y6 | Cold Play |
| 105-9318443-4200264 | Samuel R. | Illinois | 1.99 | 312 | Detroit | Y | B000002BK9 | Red Hot Chili Peppers |
| 105-1372500-0198646 | Chris G. | Massachusetts | 0.99 | 413 | New York, New York | N | B000068ZVQ | Frank Sinatra |
| 103-2628345-9238664 | Monique D. | Canada | 0.99 | 902 | Los Angeles | N | B0000010AA | Blink 182 |
| 002-1663369-6638649 | Katherine H. | Ohio | 0.99 | 440 | Beverly Hills | N | B002MXA7Q0 | Weezer |

Now we can see that these are purchase records for album download orders from Amazon. The column titles tell what has been recorded. Each row is about a particular purchase.

What information would provide a **context**? Newspaper journalists know that the lead paragraph of a good story should establish the "Five W's": *who, what, when, where,* and (if possible) *why*. Often, we add *how* to the list as well. The answers to the first two questions are essential. If we don't know *what* values are measured and *who* those values are measured on, the values are meaningless.

## Who and What

In general, the rows of a data table correspond to individual **cases** about *whom* (or about which, if they're not people) we record some characteristics. Cases go by different names, depending on the situation.

- Individuals who answer a survey are called **respondents**.
- People on whom we experiment are **subjects** or (in an attempt to acknowledge the importance of their role in the experiment) **participants**.

---

[5]www-01.ibm.com/software/data/bigdata/what-is-big-data.html

◆ Animals, plants, websites, and other inanimate subjects are often called **experimental units**.
◆ Often we simply call cases what they are: for example, *customers, economic quarters*, or *companies*.
◆ In a database, rows are called **records**—in this example, purchase records. Perhaps the most generic term is *cases*, but in any event the rows represent the *Who* of the data.

Look at all the columns to see exactly what each row refers to. Here the cases are different purchase records. You might have thought that each customer was a case, but notice that, for example, Katherine H. appears twice, in both the first and the last rows. A common place to find out exactly what each row refers to is the leftmost column. That value often identifies the cases, in this example, it's the order number. If you collect the data yourself, you'll know what the cases are. But, often, you'll be looking at data that someone else collected and you'll have to ask or figure that out yourself.

Often the cases are a **sample** from some larger **population** that we'd like to understand. Amazon doesn't care about just these customers; it wants to understand the buying patterns of *all* its customers, and, generalizing further, it wants to know how to attract other Internet users who may not have made a purchase from Amazon's site. To be able to generalize from the sample of cases to the larger population, we'll want the sample to be *representative* of that population—a kind of snapshot image of the larger world.

We must know *who* and *what* to analyze data. Without knowing these two, we don't have enough information to start. Of course, we'd always like to know more. The more we know about the data, the more we'll understand about the world. If possible, we'd like to know the *when* and *where* of data as well. Values recorded in 1803 may mean something different than similar values recorded last year. Values measured in Tanzania may differ in meaning from similar measurements made in Mexico. And knowing *why* the data were collected can tell us much about its reliability and quality.

## How the Data Are Collected

*How* the data are collected can make the difference between insight and nonsense. As we'll see later, data that come from a voluntary survey on the Internet are almost always worthless. One primary concern of statistics, to be discussed in Part III, is the design of sound methods for collecting data. Throughout this text, whenever we introduce data, we'll provide a margin note listing the W's (and H) of the data. Identifying the W's is a habit we recommend.

The first step of any data analysis is to know what you are trying to accomplish and what you want to know. To help you use statistics to understand the world and make decisions, we'll lead you through the entire process of *thinking* about the problem, *showing* what you've found, and *telling* others what you've learned. Every guided example in this text is broken into these three steps: *Think*, *Show,* and *Tell*. Identifying the problem and the *who* and *what* of the data is a key part of the *Think* step of any analysis. Make sure you know these before you proceed to *Show* or *Tell* anything about the data.

## EXAMPLE 1.1

### Identifying the *Who*

In 2015, *Consumer Reports* published an evaluation of 126 computer tablets from a variety of manufacturers.

**QUESTION:** Describe the population of interest, the sample, and the *Who* of the study.

**ANSWER:** The magazine is interested in the performance of tablets currently offered for sale. It tested a sample of 126 tablets, which are the *Who* for these data. Each tablet selected represents all similar tablets offered by that manufacturer.

# 1.3 Variables

The characteristics recorded about each individual are called **variables**. They are usually found as the columns of a data table with a name in the header that identifies what has been recorded. In the Amazon data table we find the variables *Order Number, Name, State/Country, Price*, and so on.

## Categorical Variables

Some variables just tell us what group or category each individual belongs to. Are you male or female? Pierced or not? We call variables like these **categorical**, or **qualitative**, **variables**. (You may also see them called **nominal variables** because they name categories.) Some variables are clearly categorical, like the variable *State/Country*. Its values are text and those values tell us what category the particular case falls into. But numerals are often used to label categories, so categorical variable values can also be numerals. For example, Amazon collects telephone area codes that *categorize* each phone number into a geographical region. So area code is considered a categorical variable even though it has numeric values. (But see the story in the following box.)

> *Far too many scientists have only a shaky grasp of the statistical techniques they are using. They employ them as an amateur chef employs a cookbook, believing the recipes will work without understanding why. A more cordon bleu attitude . . . might lead to fewer statistical soufflés failing to rise.*
>
> —*The Economist, June 3, 2004, "Sloppy stats shame science"*

---

### AREA CODES—NUMBERS OR CATEGORIES?

The *What* and *Why* of area codes are not as simple as they may first seem. When area codes were first introduced, AT&T was still the source of all telephone equipment, and phones had dials.

To reduce wear and tear on the dials, the area codes with the lowest digits (for which the dial would have to spin least) were assigned to the most populous regions—those with the most phone numbers and thus the area codes most likely to be dialed. New York City was assigned 212, Chicago 312, and Los Angeles 213, but rural upstate New York was given 607, Joliet was 815, and San Diego 619. For that reason, at one time the numerical value of an area code could be used to guess something about the population of its region. Since the advent of push-button phones, area codes have finally become just categories.

---

Descriptive responses to questions are often categories. For example, the responses to the questions "Who is your cell phone provider?" and "What is your marital status?" yield categorical values. When Amazon considers a special offer of free shipping to customers, it might first analyze how purchases have been shipped in the recent past. Amazon might start by counting the number of purchases shipped in each category: ground transportation, second-day air, and next-day air. Counting is a natural way to summarize a categorical variable such as *Shipping Method*. Chapter 2 discusses summaries and displays of categorical variables more fully.

## Quantitative Variables

When a variable contains measured numerical values with measurement *units*, we call it a **quantitative variable**. Quantitative variables typically record an amount or degree of something. For quantitative variables, its measurement **units** provide a meaning for the numbers. Even more important, units such as yen, cubits, carats, angstroms, nanoseconds, miles per hour, or degrees Celsius tell us the *scale* of measurement, so we know how far apart two values are. Without units, the values of a measured variable have no meaning. It does little good to be promised a raise of 5000 a year if you don't know whether it will be

paid in euros, dollars, pennies, yen, or Mauritanian Ouguiya (MRU).[6] We'll see how to display and summarize quantitative variables in Chapter 2.

Some quantitative variables don't have obvious units. The Dow Jones Industrial "Average" has units (points?) but no one talks about them. Percentages are ratios of two quantities and so the units "cancel out." But they are still percentages of something. So although it isn't imperative that a quantitative variable have explicit units, when they are not explicit, be careful to think about whether adding their values, averaging them, or otherwise treating them as numerical makes sense.

Sometimes a variable with numerical values can be treated as either categorical or quantitative depending on what we want to know from it. Amazon could record your *Age* in years. That seems quantitative, and it would be if the company wanted to know the average age of those customers who visit their site after 3 AM. But suppose Amazon wants to decide which album to feature on its site when you visit. Then thinking of your age in one of the categories Child, Teen, Adult, or Senior might be more useful. So, sometimes whether a variable is treated as categorical or quantitative is more about the question we want to ask rather than an intrinsic property of the variable itself.

## Identifiers

For a categorical variable like *Survived*, each individual is assigned one of two possible values, say *Alive* or *Dead*.[7] But for a variable with ID numbers, such as a *student ID*, each individual receives a unique value. We call a variable like this, which has exactly as many values as cases, an **identifier variable**. Identifiers are useful, but not typically for analysis.

Amazon wants to know who you are when you sign in again and doesn't want to confuse you with some other customer. So it assigns you a unique identifier. Amazon also wants to send you the right product, so it assigns a unique Amazon Standard Identification Number (ASIN) to each item it carries. You'll want to recognize when a variable is playing the role of an identifier so you aren't tempted to analyze it.

Identifier variables themselves don't tell us anything useful about their categories because we know there is exactly one individual in each. Identifiers are part of what's called **metadata**, or data about the data. Metadata are crucial in this era of large datasets because by uniquely identifying the cases, they make it possible to combine data from different sources, protect (or violate) privacy, and provide unique labels.[8] Many large databases are *relational* databases. In a relational database, different data tables link to one another by matching identifiers. In the Amazon example, the *Customer Number*, *ASIN*, and *Transaction Number* are all identifiers. The IP (Internet Protocol) address of your computer is another identifier, needed so that the electronic messages sent to you can find you.

## Ordinal Variables

A typical course evaluation survey asks, "How valuable do you think this course will be to you?" 1 = Worthless; 2 = Slightly; 3 = Middling; 4 = Reasonably; 5 = Invaluable. Is *Educational Value* categorical or quantitative? Often the best way to tell is to look to the *Why* of the study. A teacher might just count the number of students who gave each response for her course, treating *Educational Value* as a categorical variable. When she wants to see whether the course is improving, she might treat the responses as the *amount* of perceived value—in effect, treating the variable as quantitative.

But what are the units? There is certainly an *order* of perceived worth: Higher numbers indicate higher perceived worth. A course that averages 4.5 seems more valuable than one that averages 2, but we should be careful about treating *Educational Value* as purely quantitative. To treat it as quantitative, she'll have to imagine that it has "educational value units" or some similar arbitrary construct. Because there are no natural units, she

---

[6]As of 9/7/2018 $1 = 35.95 MRU.

[7]Well, maybe three values if you include Zombies.

[8]The National Security Agency (NSA) made the term "metadata" famous in 2014 by insisting that they only collected metadata on U.S. citizens' phone calls and text messages, not the calls and messages themselves. They later admitted to the bulk collection of actual data.

should be cautious. Variables that report order without natural units are often called **ordinal variables**. But saying "that's an ordinal variable" doesn't get you off the hook. You must still look to the *Why* of your study and understand what you want to learn from the variable to decide whether to treat it as categorical or quantitative.

## EXAMPLE 1.2

### Identifying the *What* and *Why* of Tablets

**RECAP:** A *Consumer Reports* article about 126 tablets lists each tablet's manufacturer, price, battery life (hrs.), the operating system (Android, iOS, or Windows), an overall quality score (0–100), and whether or not it has a memory card reader.

**QUESTION:** Are these variables categorical or quantitative? Include units where appropriate, and describe the *Why* of this investigation.

**ANSWER:** The variables are
- manufacturer (categorical)
- price (quantitative, $)
- battery life (quantitative, hrs.)
- operating system (categorical)
- quality score (quantitative, no units)
- memory card reader (categorical)

The magazine hopes to provide consumers with the information to choose a good tablet.

### JUST CHECKING

In the 2004 Tour de France, Lance Armstrong made history by winning the race for an unprecedented sixth time. In 2005, he became the only 7-time winner and set a new record for the fastest average speed—41.65 kilometers per hour. In 2012, he was banned for life for doping offenses, stripped of all of his titles and his records expunged. You can find data on all the Tour de France races in the dataset **Tour de France 2017**. Here are the first three and last seven lines of the dataset. Keep in mind that the entire dataset has over 100 entries.

1. List as many of the W's as you can for this dataset.

2. Classify each variable as categorical or quantitative; if quantitative, identify the units.

| Year | Winner | Country of Origin | Age | Team | Total Time (hours) | Avg. Speed (km/h) | Stages | Total Distance Ridden (km) | Starting Riders | Finishing Riders |
|------|--------|-------------------|-----|------|--------------------|-------------------|--------|----------------------------|-----------------|------------------|
| 1903 | Maurice Garin | France | 32 | La Française | 94.55 | 25.7 | 6 | 2428 | 60 | 21 |
| 1904 | Henri Cornet | France | 20 | Cycles JC | 96.10 | 25.3 | 6 | 2428 | 88 | 23 |
| 1905 | Louis Trousseller | France | 24 | Peugeot | 110.45 | 27.1 | 11 | 2994 | 60 | 24 |
| ... | | | | | .... | | | | | |
| 2011 | Cadel Evans | Australia | 34 | BMC | 86.21 | 39.79 | 21 | 3430 | 198 | 167 |
| 2012 | Bradley Wiggins | Great Britain | 32 | Sky | 87.58 | 39.83 | 20 | 3488 | 198 | 153 |
| 2013 | Christopher Froome | Great Britain | 28 | Sky | 94.55 | 40.55 | 21 | 3404 | 198 | 169 |
| 2014 | Vincenzo Nibali | Italy | 29 | Astana | 89.93 | 40.74 | 21 | 3663.5 | 198 | 164 |
| 2015 | Christopher Froome | Great Britain | 30 | Sky | 84.77 | 39.64 | 21 | 3660.3 | 198 | 160 |
| 2016 | Christopher Froome | Great Britain | 31 | Sky | 89.08 | 39.62 | 21 | 3529 | 198 | 174 |
| 2017 | Christopher Froome | Great Britain | 32 | Sky | 86.34 | 40.997 | 21 | 3540 | 198 | 167 |
| 2018 | Geraint Thomas | Great Britain | 32 | Sky | 83.28 | 40.210 | 21 | 3349 | 176 | 145 |

# 1.4 Models

What is a **model** for data? Models are summaries and simplifications of data that help our understanding in many ways. We'll encounter all sorts of models throughout the text. A model is a simplification of reality that gives us information that we can learn from and use, even though it doesn't represent reality exactly. A model of an airplane in a wind tunnel can give insights about the aerodynamics and flight performance of the plane even though it doesn't show every rivet.[9] In fact, it's precisely because a model is a simplification that we learn from it. Without making models for how data vary, we'd be limited to reporting only what the data we have at hand says. To have an impact on science and society we'll have to generalize those findings to the world at large.

Kepler's laws describing the motion of planets are a great example of a model for data. Using astronomical observations of Tycho Brahe, Kepler saw through the small anomalies in the measurements and came up with three simple "laws"—or models for how the planets move. Here are Brahe's observations on the declination (angle of tilt to the sun) of Mars over a twenty-year period just before 1600:

**Figure 1.1**

A plot of declination against time shows some patterns. There are many missing observations. Can you see the model that Kepler came up with from these data?



---

[9]Or tell you what movies you might see on the flight.

Here, using modern statistical methods is a plot of the model predictions from the data:

**Figure 1.2**
The model that Kepler proposed
filled in many of the missing points
and made the pattern much clearer.

**Tycho Brahe's Mars Observations**
The Orbit as Calculated with Modern Methods



Later, after Newton laid out the physics of gravity, it could be shown that the laws follow from other principles, but Kepler derived the models from data. We may not be able to come up with models as profound as Kepler's, but we'll use models throughout the text. We'll see examples of models as early as Chapter 5 and then put them to use more thoroughly later in the text when we discuss inference.

## WHAT CAN GO WRONG?

◆ **Don't label a variable as categorical or quantitative without thinking about the data and what they represent.** The same variable can sometimes take on different roles.

◆ **Don't assume that a variable is quantitative just because its values are numbers.** Categories are often given numerical labels. Don't let that fool you into thinking they have quantitative meaning. Look at the context.

◆ **Always be skeptical.** One reason to analyze data is to discover the truth. Even when you are told a context for the data, it may turn out that the truth is a bit (or even a lot) different. The context colors our interpretation of the data, so those who want to influence what you think may slant the context. A survey that seems to be about all students may in fact report just the opinions of those who visited a fan website. The question that respondents answered may be posed in a way that influences responses.

## CHAPTER REVIEW



Understand that data are values, whether numerical or labels, together with their context.

◆ *Who, what, why, where, when* (and *how*)—the W's—help nail down the context of the data.

◆ We must know *who, what,* and *why* to be able to say anything useful based on the data. The *Who* are the cases. The *What* are the variables. A variable gives information about each of the cases. The *Why* helps us decide which way to treat the variables.

◆ Stop and identify the W's whenever you have data, and be sure you can identify the cases and the variables.

Consider the source of your data and the reasons the data were collected. That can help you understand what you might be able to learn from the data.

Identify whether a variable is being used as categorical or quantitative.

- ◆ Categorical variables identify a category for each case. Usually we think about the counts of cases that fall in each category. (An exception is an identifier variable that just names each case.)
- ◆ Quantitative variables record measurements or amounts of something. They typically have units or are ratios of quantities that have units.
- ◆ Sometimes we may treat the same variable as categorical or quantitative depending on what we want to learn from it, which means some variables can't be pigeonholed as one type or the other.

## REVIEW OF TERMS

The key terms are in chapter order so you can use this list to review the material in the chapter.

| | |
|---|---|
| **Data** | Recorded values, whether numbers or labels, together with their context (p. 1). |
| **Data table** | An arrangement of data in which each row represents a case and each column represents a variable (p. 3). |
| **Context** | The context ideally tells *who* was measured, *what* was measured, *how* the data were collected, *where* the data were collected, and *when* and *why* the study was performed (p. 3). |
| **Case** | An individual about whom or which we have data (p. 3). |
| **Respondent** | Someone who answers, or responds to, a survey (p. 3). |
| **Subject** | A human experimental unit. Also called a participant (p. 3). |
| **Participant** | A human experimental unit. Also called a subject (p. 3). |
| **Experimental unit** | An individual in a study for which or for whom data values are recorded. Human experimental units are usually called subjects or participants (p. 4). |
| **Record** | Information about an individual in a database (p. 4). |
| **Sample** | A subset of a population, examined in hope of learning about the population (p. 4). |
| **Population** | The entire group of individuals or instances about whom we hope to learn (p. 4). |
| **Variable** | A variable holds information about the same characteristic for many cases (p. 5). |
| **Categorical (or qualitative) variable** | A variable that names categories with words or numerals (p. 5). |
| **Nominal variable** | The term "nominal" can be applied to a variable whose values are used only to name categories (p. 5). |
| **Quantitative variable** | A variable in which the numbers are values of measured quantities (p. 5). |
| **Unit** | A quantity or amount adopted as a standard of measurement, such as dollars, hours, or grams (p. 5). |
| **Identifier variable** | A categorical variable that records a unique value for each case, used to name or identify it (p. 6). |
| **Metadata** | Data about the data. Metadata can provide information to uniquely identify cases, making it possible to combine data from different sources, protect (or violate) privacy, and label cases uniquely (p. 6). |
| **Ordinal variable** | The term "ordinal" can be applied to a variable whose categorical values possess some kind of order (p. 7). |
| **Model** | A description or representation, in mathematical and statistical terms, of the behavior of a phenomenon based on data (p. 8). |

## TECH SUPPORT

## Entering Data

These days, nobody does statistics by hand. We use technology: a programmable calculator or a statistics program on a computer. Professionals all use a *statistics package* designed for the purpose. We will provide many examples of results from a statistics package throughout the text. Rather than choosing one in particular, we'll offer generic results that look like those produced by all the major statistics packages but don't exactly match any of them. Then, in the Tech Support section at the end of each chapter, we'll provide hints for getting started on several of the major packages.

If you understand what the computer needs to know to do what you want and what it needs to show you in return, you can figure out the specific details of most packages pretty easily.

For example, to get your data into a computer statistics package, you need to tell the computer:

▶ Where to find the data. This usually means directing the computer to a file stored on your computer's disk or to data on a database. Or it might just mean that you have copied the data from a spreadsheet program or Internet site and it is currently on your computer's clipboard. Usually, the data should be in the form of a data table with cases in the rows and variables in the columns. Most computer statistics packages prefer the *delimiter* that marks the division between elements of a data table to be a tab character (comma is another common delimiter) and the delimiter that marks the end of a case to be a *return* character. The data used in this text can be found in the DASL archive at dasl.datadescription.com and on the text's website at media.pearsoncmg.com/aw/aw_deveaux_stats_5/cw/statdm5d_home.html.

▶ Where to put the data. (Usually this is handled automatically.)

▶ What to call the variables. Some data tables have variable names as the first row of the data, and often statistics packages can take the variable names from the first row automatically.

▶ Excel is often used to help organize, manipulate, and prepare data for other software packages. Many of the other packages take Excel files as inputs. Alternatively, you can copy a data table from Excel and paste it into many packages, or export Excel spreadsheets as tab delimited (.txt) or comma delimited files (.csv), which can be easily shared and imported into other programs. All data files provided with this text are in tab-delimited text (.txt) format.

### EXCEL

To open a file containing data in Excel:

▶ Choose **File** > **Open**.

▶ Browse to find the file to open. Excel supports many file formats.

▶ Other programs can import data from a variety of file formats, but all can read both tab delimited (.txt) and comma delimited (.csv) text files.

▶ You can also copy tables of data from other sources, such as Internet sites, and paste them into an Excel spreadsheet. Excel can recognize the format of many tables copied this way, but this method may not work for some tables.

▶ Excel may not recognize the format of the data. If data include dates or other special formats ($, €, ¥, etc.), identify the desired format. Select the cells or columns to reformat and choose **Format** > **Cell**. Often, the General format is the best option for data you plan to move to a statistics package.

### DATA DESK

To read data into Data Desk:

▶ Click the **Open File** icon or choose **File** > **Open**. The dialog lets you specify variable names (or take them from the first row of the data), the delimiter, or how to read formatted data.

▶ **File** > **Import** works the same way, but instead of starting a new data file, it adds the data in the file to the current data file. Data Desk can work with multiple data tables in the same file.

▶ If the data are already in another program, such as, for example, a spreadsheet, **Copy** the data table (including the column headings). In Data Desk choose **Edit** > **Paste** variables. There is no need to create variables first; Data Desk does that automatically. You'll see the same dialog as for Open and Import.

▶ The DASL library of datasets (dasl.datadescription.com) provides direct links to transfer data to Data Desk.

## JMP

To import a text file:

▶ Choose **File** > **Open** and select the file from the dialog. At the bottom of the dialog screen you'll see **Open As:**—be sure to change to **Data (Using Preview)**. This will allow you to specify the delimiter and make sure the variable names are correct. (**JMP** also allows various formats to be imported directly, including .xls files.)

You can also paste a dataset in directly (with or without variable names) by selecting:

▶ **File** > **New** > **New Data Table** and then **Edit** > **Paste** (or **Paste with Column Names** if you copied the names of the variables as well).

Finally, you can import a dataset from a URL directly by selecting:

▶ **File** > **Internet Open** and pasting in the address of the website. JMP will attempt to find data on the page. It may take a few tries and some edits to get the dataset in correctly.

## MINITAB

To import a text or Excel file:

▶ Choose **File** > **Open Worksheet**. From **Files of type**, choose **Text (*.txt)** or **Excel (*.xls; *xlsx)**.

▶ Browse to find and select the file.

▶ In the lower right corner of the dialog, choose **Open** to open the data file alone, or **Merge** to add the data to an existing worksheet.

▶ Click **Open**.

## R

**R** can import many types of files, but text files (tab or comma delimited) are easiest. If the file is tab delimited and contains the variable names in the first row, then:

> **mydata = read.delim(file.choose())**

will give a dialog where you can pick the file you want to import. It will then be in a data frame called mydata. If the file is comma delimited, use:

> **mydata = read.csv(file.choose())**

**COMMENTS**

RStudio provides an interactive dialog that may be easier to use. For other options, including the case that the file does not contain variable names, consult **R** help.

## SPSS

To import a text file:

▶ Choose **File** > **Open** > **Data**. Under "Files of type," choose **Text (*.txt,*.dat)**. Select the file you want to import. Click **Open**.

▶ A window will open called **Text Import Wizard**. Follow the steps, depending on the type of file you want to import.

## STATCRUNCH

StatCrunch offers several ways to enter data. Click **MyStatCrunch** > **My Data**. Click a dataset to analyze the data or edit its properties.

Click a dataset link to analyze the data or edit its properties to import a new dataset.

▶ Choose **Select a file on my computer**,

▶ Enter the URL of a file,

▶ Paste data into a form, or

▶ Type or paste data into a blank data table.

For the "select a file on my computer" option, StatCrunch offers a choice of space, comma, tab, or semicolon delimiters. You may also choose to use the first line as the names of the variables.

After making your choices, select the **Load File** button at the bottom of the screen.

StatCrunch has direct access to the datasets on the text's website.

# EXERCISES

## SECTION 1.1

**1. Grocery shopping**  Many grocery store chains offer customers a card they can scan when they check out and offer discounts to people who do so. To get the card, customers must give information, including a mailing address and e-mail address. The actual purpose is not to reward loyal customers but to gather data. What data do these cards allow stores to gather, and why would they want that data?

**2. Online shopping**  Online retailers such as Amazon.com keep data on products that customers buy, and even products they look at. What does Amazon hope to gain from such information?

**3. Parking lots**  Sensors in parking lots are able to detect and communicate when spaces are filled in a large covered parking garage next to an urban shopping mall. How might the owners of the parking garage use this information both to attract customers and to help the store owners in the mall make business plans?

**4. Satellites and global climate change**  Satellites send back nearly continuous data on the earth's land masses, oceans, and atmosphere from space. How might researchers use this information in both the short and long terms to help study changes in the earth's climate?

## SECTION 1.2

**5. Super Bowl**  Sports announcers love to quote statistics. During the Super Bowl, they particularly love to announce when a record has been broken. They might have a list of all Super Bowl games, along with the scores of each team, total scores for the two teams, margin of victory, passing yards for the quarterbacks, and many more bits of information. Identify the *Who* in this list.

**6. Nobel laureates**  The website www.nobelprize.org allows you to look up all the Nobel prizes awarded in any year. The data are not listed in a table. Rather you drag a slider to the year and see a list of the awardees for that year. Describe the *Who* in this scenario.

**7. Health records**  The National Center for Health Statistics (NCHS) conducts an extensive survey consisting of an interview and medical examination with a representative sample of about 5000 people a year. The interview includes demographic, socioeconomic, dietary, and other health-related questions. The examination "consists of medical, dental, and physiological measurements, as well as laboratory tests administered by highly trained medical personnel" (www.cdc.gov/nchs/nhanes/about_nhanes.htm). Describe the sample, the population, the *Who* and the *What* of this study.

**8. Facebook.**  Facebook uploads more than 350 million photos every day onto its servers. For this collection, describe the *Who* and the *What*.

## SECTION 1.3

**9. Grade levels**  A person's grade in school is generally identified by a number.

  a) Give an example of a *Why* in which grade level is treated as categorical.
  b) Give an example of a *Why* in which grade level is treated as quantitative.

**10. ZIP codes**  The U.S. Postal Service uses five-digit ZIP codes to identify locations to assist in delivering mail.

  a) In what sense are ZIP codes categorical?
  b) Is there any ordinal sense to ZIP codes? In other words, does a higher ZIP code tell you anything about a location compared to a lower ZIP code?

**11. Voters**  A February 2010 Gallup Poll question asked, "In politics, as of today, do you consider yourself a Republican, a Democrat, or an Independent?" The possible responses were "Democrat," "Republican," "Independent," "Other," and "No Response." What kind of variable is the response?

**12. Job hunting**  A June 2011 Gallup Poll asked Americans, "Thinking about the job situation in America today, would you say that it is now a good time or a bad time to find a quality job?" The choices were "Good time" or "Bad time." What kind of variable is the response?

**13. Medicine**  A pharmaceutical company conducts an experiment in which a subject takes 100 mg of a substance orally. The researchers measure how many minutes it takes for half of the substance to exit the bloodstream. What kind of variable is the company studying?

**14. Stress**  A medical researcher measures the increase in heart rate of patients who are taking a stress test. What kind of variable is the researcher studying?

## SECTION 1.4

**15. Voting and elections**  Pollsters are interested in predicting the outcome of elections. Give an example of how they might model whether someone is likely to vote.

**16. Weather**  Meteorologists utilize sophisticated models to predict the weather up to ten days in advance. Give an example of how they might assess their models.

**17. The news**  Find a newspaper or magazine article in which some data are reported. For the data discussed in the article, identify as many of the W's as you can. Include a copy of the article with your report.

**18. The Internet**  Find an Internet source that reports on a study and describes the data. Print out the description and identify as many of the W's as you can.

*(Exercises 19–26) For each description of data, identify Who and What were investigated and the population of interest.*

**19. Gaydar** A study conducted by a team of American and Canadian researchers found that during ovulation, a woman can tell whether a man is gay or straight by looking at his face. To explore the subject, the authors conducted three investigations, the first of which involved 40 undergraduate women who were asked to guess the sexual orientation of 80 men based on photos of their face. Half of the men were gay, and the other half were straight. All held similar expressions in the photos or were deemed to be equally attractive. None of the women were using any contraceptive drugs at the time of the test. The result: the closer a woman was to her peak ovulation, the more accurate her guess. (health.usnews.com/health-news/family-health/brain-and-behavior/articles/2011/06/27/ovulation-seems-to-aid-womens-gaydar)

**20. Hula-hoops** The hula-hoop, a popular children's toy in the 1950s, has gained popularity as an exercise in recent years. But does it work? To answer this question, the American Council on Exercise conducted a study to evaluate the cardio and calorie-burning benefits of "hooping." Researchers recorded heart rate and oxygen consumption of participants, as well as their individual ratings of perceived exertion, at regular intervals during a 30-minute workout. (www.acefitness.org/certifiednewsarticle/1094/)

**21. Bicycle safety** Ian Walker, a psychologist at the University of Bath, wondered whether drivers treat bicycle riders differently when they wear helmets. He rigged his bicycle with an ultrasonic sensor that could measure how close each car was that passed him. He then rode on alternating days with and without a helmet. Out of 2500 cars passing him, he found that when he wore his helmet, motorists passed 3.35 inches closer to him, on average, than when his head was bare. (Source: *NY Times*, Dec. 10, 2006)

**22. Investments** Some companies offer 401(k) retirement plans to employees, permitting them to shift part of their before-tax salaries into investments such as mutual funds. Employers typically match 50% of the employees' contribution up to about 6% of salary. One company, concerned with what it believed was a low employee participation rate in its 401(k) plan, sampled 30 other companies with similar plans and asked for their 401(k) participation rates.

**23. Honesty** Coffee stations in offices often just ask users to leave money in a tray to pay for their coffee, but many people cheat. Researchers at Newcastle University alternately taped two posters over the coffee station. During one week, it was a picture of flowers; during the other, it was a pair of staring eyes. They found that the average contribution was significantly higher when the eyes poster was up than when the flowers were there. Apparently, the mere feeling of being watched—even by eyes that were not real—was enough to encourage people to behave more honestly. (Source: *NY Times*, Dec. 10, 2006)

**24. Blindness** A study begun in 2011 examines the use of stem cells in treating two forms of blindness, Stargardt's disease and dry age-related macular degeneration. Each of the 24 patients entered one of two separate trials in which embryonic stem cells were to be used to treat the condition. (www.blindness.org/index.php?view=article&id=2514:stem-cell-clinical-trial-for-stargardt-disease-set-to-begin-&option=com_content&Itemid=122)

**25. Not-so-diet soda** A look at 474 participants in the San Antonio Longitudinal Study of Aging found that participants who drank two or more diet sodas a day "experienced waist size increases six times greater than those of people who didn't drink diet soda." (*J Am Geriatr Soc.* 2015 Apr;63(4):708–15. doi: 10.1111/jgs.13376. Epub 2015 Mar 17.)

**26. Molten iron** The Cleveland Casting Plant is a large, highly automated producer of gray and nodular iron automotive castings for Ford Motor Company. The company is interested in keeping the pouring temperature of the molten iron (in degrees Fahrenheit) close to the specified value of 2550 degrees. Cleveland Casting measured the pouring temperature for 10 randomly selected crankshafts.

*(Exercises 27–40) For each description of data, identify the W's, name the variables, specify for each variable whether its use indicates that it should be treated as categorical or quantitative, and, for any quantitative variable, identify the units in which it was measured (or note that they were not provided).*

**27. Weighing bears** Because of the difficulty of weighing a bear in the woods, researchers caught and measured 54 bears, recording their weight, neck size, length, and sex. They hoped to find a way to estimate weight from the other, more easily determined quantities.

**28. Schools** The State Education Department requires local school districts to keep these records on all students: age, race or ethnicity, days absent, current grade level, standardized test scores in reading and mathematics, and any disabilities or special educational needs.

**29. Arby's menu** A listing posted by the Arby's restaurant chain gives, for each of the sandwiches it sells, the type of meat in the sandwich, the number of calories, and the serving size in ounces. The data might be used to assess the nutritional value of the different sandwiches.

**30. Age and party** The Gallup Poll conducted a representative telephone survey of 1180 American voters during the first quarter of 2007. Among the reported results were the voter's region (Northeast, South, etc.), age, party affiliation, and whether or not the person had voted in the 2006 midterm congressional election.

**31. Babies** Medical researchers at a large city hospital investigating the impact of prenatal care on newborn health collected data from 882 births during 1998–2000. They kept track of the mother's age, the number of weeks the pregnancy lasted, the type of birth (cesarean, induced, natural), the level of prenatal care the mother had (none, minimal, adequate), the birth weight and sex of the baby, and whether the baby exhibited health problems (none, minor, major).

**32. Flowers** In a study appearing in the journal *Science,* a research team reports that plants in southern England are flowering earlier in the spring. Records of the first flowering dates for 385 species over a period of 47 years show that flowering has advanced an average of 15 days per decade, an indication of climate warming, according to the authors.

**33. Herbal medicine** Scientists at a major pharmaceutical firm conducted an experiment to study the effectiveness of an herbal compound to treat the common cold. They exposed each patient to a cold virus, then gave them either the herbal compound or a sugar solution known to have no effect on colds. Several days

| Year | Winner | Jockey | Trainer | Owner | Time |
|------|--------|--------|---------|-------|------|
| 1875 | Aristides | O. Lewis | A. Williams | H. P. McGrath | 2:37.75 |
| 1876 | Vagrant | R. Swim | J. Williams | William Astor | 2:38.25 |
| 1877 | Baden Baden | W. Walker | E. Brown | Daniel Swigert | 2:38 |
| 1878 | Day Star | J. Carter | L. Paul | T. J. Nichols | 2:37.25 |
| … | | | | | |
| 2011 | Animal Kingdom | J. Velazquez | H. G. Motion | Team Valor | 2:02.04 |
| 2012 | I'll Have Another | M. Gutierrez | D. O'Neill | Reddam Racing | 2:01.83 |
| 2013 | Orb | J. Rosario | S. McGaughey | Stuart Janney & Phipps Stable | 2:02.89 |
| 2014 | California Chrome | Victor Espinoza | Art Sherman | California Chrome, LLC | 2:03.66 |
| 2015 | American Pharoah | Victor Espinoza | Bob Baffert | Zayat Stables, LLC | 2:03.03 |
| 2016 | Nyquist | M. Gutierrez | Doug F. O'Neill | Reddam Racing LLC | 2:01.31 |
| 2017 | Always Dreaming | J. Velazquez | Todd Pletcher | Meb Racing Stables | 2:03.59 |
| 2018 | Justify | M. Smith | Bob Baffert | China Horse Club | 2:04.20 |

*Source:* Excerpt from HorseHats.com. Published by Thoroughbred Promotions.

later they assessed each patient's condition, using a cold severity scale ranging from 0 to 5. They found no evidence of benefits of the compound.

**34. Vineyards**  Business analysts hoping to provide information helpful to American grape growers compiled these data about vineyards: size (acres), number of years in existence, state, varieties of grapes grown, average case price, gross sales, and percent profit.

**35. Streams**  In performing research for an ecology class, students at a college in upstate New York collect data on streams each year. They record a number of biological, chemical, and physical variables, including the stream name, the substrate of the stream (limestone, shale, or mixed), the acidity of the water (pH), the temperature (°C), and the BCI (a numerical measure of biological diversity).

**36. Fuel economy**  The Environmental Protection Agency (EPA) tracks fuel economy of automobiles based on information from the manufacturers (Ford, Toyota, etc.). Among the data the agency collects are the manufacturer, vehicle type (car, SUV, etc.), weight, horsepower, and gas mileage (mpg) for city and highway driving.

**37. Refrigerators**  In 2013, *Consumer Reports* published an article evaluating refrigerators. It listed 353 models, giving the brand, cost, size (cu ft), type (such as top freezer), estimated annual energy cost, an overall rating (good, excellent, etc.), and the repair history for that brand (percentage requiring repairs over the past 5 years).

**38. Walking in circles**  People who get lost in the desert, mountains, or woods often seem to wander in circles rather than walk in straight lines. To see whether people naturally walk in circles in the absence of visual clues, researcher Andrea Axtell tested 32 people on a football field. One at a time, they stood at the center of one goal line, were blindfolded, and then tried to walk to the other goal line. She recorded each individual's sex, height, handedness, the number of yards each was able to walk before going out of bounds, and whether each wandered off course to the left or right. No one made it all the way to the far end of the field without crossing one of the sidelines. (Source: *STATS* No. 39, Winter 2004)

**T 39. Kentucky Derby 2018**  The Kentucky Derby is a horse race that has been run every year since 1875 at Churchill Downs in Louisville, Kentucky. The race started as a 1.5-mile race, but in 1896, it was shortened to 1.25 miles because experts felt that 3-year-old horses shouldn't run such a long race that early in the season. (It has been run in May every year but one—1901—when it took place on April 29.) Above are the data for the first four and eight recent races.

**T 40. Indy 500 2018**  The 2.5-mile Indianapolis Motor Speedway has been the home to a race on Memorial Day nearly every year since 1911. Even during the first race, there were controversies. Ralph Mulford was given the checkered flag first but took three extra laps just to make sure he'd completed 500 miles. When he finished, another driver, Ray Harroun, was being presented with the winner's trophy, and Mulford's protests were ignored. Harroun averaged 74.6 mph for the 500 miles. In 2013, the winner, Tony Kanaan, averaged over 187 mph, beating the previous record by over 17 mph!

Here are the data for the first five races and six recent Indianapolis 500 races.

| Year | Driver | Time (hr:min:sec) | Speed (mph) |
|------|--------|-------------------|-------------|
| 1911 | Ray Harroun | 6:42:08 | 74.602 |
| 1912 | Joe Dawson | 6:21:06 | 78.719 |
| 1913 | Jules Goux | 6:35:05 | 75.933 |
| 1914 | René Thomas | 6:03:45 | 82.474 |
| 1915 | Ralph DePalma | 5:33:55.51 | 89.840 |
| … | | | |
| 2013 | Tony Kanaan | 2:40:03.4181 | 187.433 |
| 2014 | Ryan Hunter-Reay | 2:40:48.2305 | 186.563 |
| 2015 | Juan Pablo Montoya | 3:05:56.5286 | 161.341 |
| 2016 | Alexander Rossi | 3:00:02.0872 | 166.634 |
| 2017 | Takuma Sato | 3:13:3.3584 | 155.395 |
| 2018 | Will Power | 2:59:42.6365 | 166.935 |

**T** **41. Kentucky Derby 2018 on the computer** Load the **Kentucky Derby 2018** data into your preferred statistics package and answer the following questions;

    a) What was the name of the winning horse in 1880?
    b) When did the length of the race change?
    c) What was the winning time in 1974?
    d) Only one horse has run the Derby in less than 2 minutes. Which horse and in what year?

**T** **42. Indy 500 2018 on the computer** Load the **Indy 500 2018** data into your preferred statistics package and answer the following questions:

    a) What was the average speed of the winner in 1920?
    b) How many times did Bill Vukovich win the race in the 1950s?
    c) How many races took place during the 1940s?

## JUST CHECKING

**Answers**

1. *Who*—Tour de France races; *What*—year, winner, country of origin, age, team, total time, average speed, stages, total distance ridden, starting riders, finishing riders; *How*—official statistics at race; *Where*—France (for the most part); *When*—1903 to 2016; *Why*—not specified (To see progress in speeds of cycling racing?)

2.

| Variable | Type | Units |
|---|---|---|
| Year | Quantitative or Identifier | Years |
| Winner | Categorical | |
| Country of Origin | Categorical | |
| Age | Quantitative | Years |
| Team | Categorical | |
| Total Time | Quantitative | Hours/minutes/seconds |
| Average Speed | Quantitative | Kilometers per hour |
| Stages | Quantitative | Counts (stages) |
| Total Distance | Quantitative | Kilometers |
| Starting Riders | Quantitative | Counts (riders) |
| Finishing Riders | Quantitative | Counts (riders) |

# 2

# Displaying and Describing Data

## WHERE ARE WE GOING?

We can summarize and describe data values in a variety of ways. You'll probably recognize these displays and summaries. This chapter is a fast review of these concepts so we all agree on terms, notation, and methods. We'll be using these displays and descriptions throughout the rest of the text.

What happened on the *Titanic* at 11:40 on the night of April 14, 1912, is well known. Frederick Fleet's cry of "Iceberg, right ahead" and the three accompanying pulls of the crow's nest bell signaled the beginning of a nightmare that has become legend. By 2:15 AM, the *Titanic*, thought by many to be unsinkable, had sunk. Only 712 of the 2208 people on board survived. The others (nearly 1500) met their icy fate in the cold waters of the North Atlantic.

Table 2.1 shows some data about the passengers and crew aboard the *Titanic*. Each case (row) of the data table represents a person on board the ship. The variables are the person's *Name, Survival* status (Dead or Alive), *Age* (in years), *Age Category* (Adult or Child), *Sex* (Male or Female), *Price* Paid (in British pounds, **£**), and ticket *Class* (First, Second, Third, or Crew). Some of these, such as *Age* and *Price*, record numbers. These are called

**Table 2.1**

Part of a data table showing seven variables for 11 people aboard the *Titanic*.

| Name | Survived | Age | Adult/Child | Sex | Price (£) | Class |
|------|----------|-----|-------------|-----|-----------|-------|
| ABBING, Mr Anthony | Dead | 42 | Adult | Male | 7.55 | 3 |
| ABBOTT, Mr Ernest Owen | Dead | 21 | Adult | Male | 0 | Crew |
| ABBOTT, Mr Eugene Joseph | Dead | 14 | Child | Male | 20.25 | 3 |
| ABBOTT, Mr Rossmore Edward | Dead | 16 | Adult | Male | 20.25 | 3 |
| ABBOTT, Mrs Rhoda Mary "Rosa" | Alive | 39 | Adult | Female | 20.25 | 3 |
| ABELSETH, Miss Karen Marie | Alive | 16 | Adult | Female | 7.65 | 3 |
| ABELSETH, Mr Olaus Jörgensen | Alive | 25 | Adult | Male | 7.65 | 3 |
| ABELSON, Mr Samuel | Dead | 30 | Adult | Male | 24 | 2 |
| ABELSON, Mrs Hannah | Alive | 28 | Adult | Female | 24 | 2 |
| ABRAHAMSSON, Mr Abraham August Johannes | Alive | 20 | Adult | Male | 7.93 | 3 |
| ABRAHIM, Mrs Mary Sophie Halaut | Alive | 18 | Adult | Female | 7.23 | 3 |

**quantitative** variables. Others, like *Survival* and *Class*, place each case in a single category, and are called **categorical** variables. (Data in **Titanic**)

The problem with a data table like this—and in fact with all data tables—is that you can't *see* what's going on. And seeing is just what we want to do. We need ways to show the data so that we can see patterns, relationships, trends, and exceptions.

## The Three Rules of Data Analysis

There are three things you should always do first with data:

1. **Make a picture.** A display of your data will reveal things you're not likely to see in a table of numbers and will help you *Think* clearly about the patterns and relationships that may be hiding in your data.
2. **Make a picture.** A well-designed display will *Show* the important features and patterns in your data. It could also show you things you did not expect to see: extraordinary (possibly wrong) data values or unexpected patterns.
3. **Make a picture.** The best way to *Tell* others about your data is with a well-chosen picture.

These are the three rules of data analysis. There are pictures of data throughout the text, and new kinds keep showing up. These days, technology makes drawing pictures of data easy, so there is no reason not to follow the three rules.

We make graphs for two primary reasons: to understand more about data and to show others what we have learned and want them to understand. The first reason calls for simple graphs with little adornment; the second often uses visually appealing additions to draw the viewer's attention. Regardless of their function, graphs should be easy to read and understand and should represent the facts of the data honestly. Axes should be clearly labeled with the names of the variables they display. The intervals set off by "tick marks" should occur at values easy to think about: 5, 10, 15, and 20 are simpler marks than, say, 1.7, 2.3, 2.9, and 3.5. And tick labels that run for several digits are almost never a good idea. Graphs should have a "key" that identifies colors and symbols if those are meaningful in the graph. And all graphs should carry a title or caption that says what the graph displays and suggests what about it is salient or important.

## The Area Principle

A bad picture can distort our understanding rather than help it. What impression do you get from Figure 2.1 about who was aboard the ship?

| WHO | People on the *Titanic* |
|-----|------|
| WHAT | Name, survival status, age, adult/child, sex, price paid, ticket class |
| WHEN | April 14, 1912 |
| WHERE | North Atlantic |
| HOW | www.encyclopedia-titanica.org |
| WHY | Historical interest |

**Figure 2.1**

How many people were in each class on the *Titanic*? From this display, it looks as though the service must have been great, since most aboard were crew members. Although the length of each ship here corresponds to the correct number, the impression is all wrong. In fact, only about 40% were crew.

The *Titanic* was certainly a luxurious ship, especially for those in first class, but Figure 2.1 gives the mistaken impression that most of the people on the *Titanic* were crew members, with a few passengers along for the ride. What's wrong? The lengths of the ships *do* match the number of people in each ticket class category. However, our eyes tend to be more impressed by the *area* than by other aspects of each ship image. So, even though the *length* of each ship matches up with one of the totals, it's the associated *area* in the image that we notice. There were about 3 times as many crew as second-class passengers, and the ship depicting the number of crew members is about 3 times longer than the ship depicting second-class passengers. The problem is that it occupies about 9 times the area. That just isn't a correct impression.

The best data displays observe a fundamental principle of graphing data called the **area principle**. The area principle says that the area occupied by a part of the graph should correspond to the magnitude of the value it represents. Violations of the area principle are a common way to lie (or, since most mistakes are unintentional, we should say err) with statistics.

# 2.1  Summarizing and Displaying a Categorical Variable

## Frequency Tables

Categorical variables are easy to summarize in a **frequency table** that lists the number of cases in each category along with its name.

For ticket *Class*, the categories are First, Second, Third, and Crew:

| Class | Count |
|-------|-------|
| First | 324 |
| Second | 285 |
| Third | 710 |
| Crew | 889 |

**Table 2.2**
A frequency table of the *Titanic* passengers.

| Class | Percentage (%) |
|-------|----------------|
| First | 14.67 |
| Second | 12.91 |
| Third | 32.16 |
| Crew | 40.26 |

**Table 2.3**
A relative frequency table for the same data.

A **relative frequency table** displays *percentages* (or *proportions*) rather than the counts in each category. Both types of tables show the **distribution** of a categorical variable because they name the possible categories and tell how frequently each occurs. (The percentages should total 100%, although the sum may be a bit too high or low if the individual category percentages have been rounded.)

## Bar Charts

Although not as visually entertaining as the ships in Figure 2.1, the **bar chart** in Figure 2.2 gives an *accurate* visual impression of the distribution because it obeys the area principle. Now it's easy to see that the majority of people on board were *not* crew. We can also see that there were about 3 times as many crew members as second-class passengers. And there were more than twice as many third-class passengers as either first- or second-class passengers—something you may have missed in the frequency table. Bar charts make these kinds of comparisons easy and natural.



**Figure 2.2**
*People on the* Titanic *by Ticket Class.* With the area principle satisfied, we can see the true distribution more clearly.

## EXAMPLE 2.1

### What Do You Think of Congress?

In December 2017, the Gallup survey asked 1049 people how they viewed a variety of professions. Specifically they asked, "How would you rate the honesty and ethical standards of people in these different fields?" For Members of Congress, the results were

| Rating | Percentage (%) |
|---|---|
| Very high | 2 |
| High | 9 |
| Average | 29 |
| Low | 36 |
| Very low | 24 |
| No opinion | 1 |

**QUESTION:** What kind of table is this? What would be an appropriate display?

**ANSWER:** This is a relative frequency table because the numbers displayed are percentages, not counts. A bar chart would be appropriate:

**Opinions About Members of Congress**



## EXAMPLE 2.2

### Which Gadgets Do You Use?

In 2014, the Pew Research Organization asked 1005 U.S. adults which of the following electronic items they use: cell phone, smartphone, computer, handheld e-book reader (e.g., Kindle or Nook), or tablet. The results were

| Device | Percentage (%) using the device |
|---|---|
| Cell phone | 86.8 |
| Smartphone | 54.0 |
| Computer | 77.5 |
| E-book reader | 32.2 |
| Tablet | 41.9 |

**QUESTION:** Is this a frequency table, a relative frequency table, or neither? How could you display these data?

**ANSWER:** This is not a frequency table because the numbers displayed are not counts. Although the numbers are percentages, they do not sum to 100%. A person can use more than one device, so this is not a relative frequency table either. A bar chart might still be appropriate, but the numbers do not sum to 100%.

**Percentage Using Each Device**

## Pie Charts

**Pie charts** display all the cases as a circle whose slices have areas proportional to each category's fraction of the whole.

Pie charts give a quick impression of the distribution. Because we're used to cutting up pies into 2, 4, or 8 pieces, pie charts are particularly good for seeing relative frequencies near 1/2, 1/4, or 1/8.

Bar charts are almost always better than pie charts for comparing the relative frequencies of categories. Pie charts are widely understood and colorful, and they often appear in reports, but Figure 2.3 shows why statisticians prefer bar charts.

**Figure 2.3**
Pie charts may be attractive, but it can be hard to see patterns in them. Can you discern the differences in distributions depicted by these pie charts?



**Figure 2.4**
Bar charts of the same values as shown in Figure 2.3 make it much easier to compare frequencies in groups.

## Ring Charts

A ring (or donut) chart is a modified form of pie chart that displays only the "crust" of the pie—a ring that is partitioned into regions proportional in area to each value. You can think of the ring as the bars of a bar chart stuck end to end and wrapped around the circle. Ring charts are somewhere between bar charts and pie charts. They may be easier to read (or not). Judge for yourself:



**Figure 2.5**
Ring charts compromise between pie and bar charts. These ring charts show the same values as the pie charts in Figure 2.3. Do you find it easier to see the patterns?

### RANDOM MATTERS

Is it random, or is something systematic going on? Separating the *signal* (the systematic) from the *noise* (the random) is a fundamental skill of statistics.

A geoscientist notices that global temperatures have increased steadily during the past 50 years. Could the pattern be random, or is the earth warming?

An analyst notices that the stock market seems to go up more often on Tuesday afternoons when it rains in Chicago. Is that something she should bank on?

One of the challenges to answering questions like these is that we have only one earth and one stock market history. What if we had two? Or many? Sometimes we can use a computer to *simulate* other situations, to pretend that we have more than one realization of a phenomenon. In these *Random Matters* sections, we'll use the computer as our lab to test what might happen if we could repeat our data collection many times.

You probably know that the "rules of the sea" were enforced on the *Titanic*—women and children were allowed to board the *Titanic* lifeboats before the men. Did ticket class (first, second, or third) also make a difference? Suppose the 712 survivors were chosen at random, giving everyone an equal chance to get into a lifeboat. Would the distribution have been different? Let's look. We selected 712 people at random from the list of those aboard the *Titanic* and made a pie chart of ticket class. We repeated the random selection 24 times, making a new pie chart of each selected group of 712 passengers. Among these pie charts in Figure 2.6 we've "hidden" the actual distribution of survivors. Can you pick out the real distribution? If so, then that might convince you that the lifeboats weren't filled randomly, with everyone getting an equal chance.

**Figure 2.6**
The distribution of ticket class in 24 simulated lifeboats and the actual distribution of survivors. Can you find the real one? If so, this suggests that people didn't all have an equal chance to survive.[1]

In this example, the difference is pretty obvious. There were more survivors from first and second class and fewer from third class than there would have been were everyone given an equal chance. In other situations, the differences may not be as obvious, so we'll need to develop more sophisticated tools to help distinguish signals from noise.

# 2.2 Displaying a Quantitative Variable

## Histograms

How can we make a bar chart for a quantitative variable? We can't, because quantitative variables don't have categories. Instead, we make a **histogram**.

Histograms and bar charts both use bars, but they are fundamentally different. The bars of a bar chart display the count for each category, so they could be arranged in any order[2]

---

[1]Wait. Didn't we just say we prefer bar charts? Well, sometimes pie charts are actually a good choice. They are compact, colorful, and—most important—they satisfy the area principle. A figure with 25 bar charts would look much more confusing.
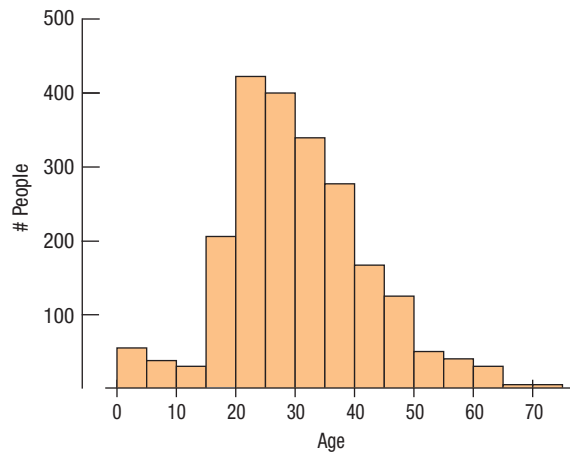
[2]Many statistics programs choose alphabetical order, which is rarely the most useful one.

(and should be displayed with a space between them). The horizontal axis of a bar chart just names the categories. The horizontal axis of a histogram shows the values of the variable in order. A histogram slices up that axis into equal-width bins, and the bars show the counts for each bin. Now **gaps** are meaningful; they show regions with no observations.

Figure 2.7 shows a histogram of the ages of those aboard the *Titanic*. In this histogram, each bin has a width of 5 years, so, for example, the height of the tallest bar shows that the most populous age group was the 20- to 24-year-olds, with over 400 people.[3] The youngest passengers were infants, and the oldest was more than 70 years old.

**Figure 2.7**
A histogram of the distribution of ages of those aboard *Titanic*.



The fact that there are fewer and fewer people in the 5-year bins from age 25 on probably doesn't surprise you either. After all, there are increasingly fewer people of advancing age in the general population as well, and there were no very elderly people on board the *Titanic*. But the bins on the left are a little strange. It looks like there were more infants and toddlers (0–5 years old) than there were preteens.

Does this distribution look plausible? You may not have guessed that fact about the infants and preteens, but it doesn't seem out of the question. It is often a good idea to imagine what the distribution might look like before you make the display. That way you'll be less likely to be fooled by errors in the data or when you accidentally graph the wrong variable.[4]

<div style="border-left: 4px solid green;">

**EXAMPLE 2.3**

## Earthquakes and Tsunamis

In 2011, the most powerful earthquake ever recorded in Japan created a wall of water that devastated the northeast coast of Japan and left nearly 25,000 people dead or missing. The 2011 tsunami in Japan was caused by a 9.0 magnitude earthquake. It was particularly noted for the damage it caused to the Fukushima Daiichi nuclear power plant, causing a core meltdown and international concern. As disastrous as it was, the Japan tsunami was not nearly as deadly as the 2004 tsunami on the west coast of Sumatra that killed an estimated 227,899 people, making it the

</div>

---

[3]The histogram bar appears to go from 20 to 25, but most statistics programs include values at the lower limit in the bar and put values at the upper limit in the next bin.

[4]You'll notice that we didn't say *if* you graph the wrong variable, but rather *when*. Everyone makes mistakes, and you'll make your share. But if you always think about what your graph or analysis says about the world and judge whether that is reasonable, you can catch many errors before they get away.

most lethal tsunami on record. The earthquake that caused it had magnitude 9.1—more than 25% more powerful than the Japanese earthquake. Were these earthquakes truly extraordinary, or did they just happen at unlucky times and places? The magnitudes (measured or estimated) are available for 999 of the 1116 earthquakes known to have caused tsunamis, dating back to 426 BCE. (Data in **Tsunamis 2018**)

**QUESTION:** What can we learn from these data?

| WHO | 1116 earthquakes known or suspected to have caused tsunamis for which we have data or good estimates |
| --- | --- |
| WHAT | Magnitude (Richter scale), depth (m), date, location, and other variables |
| WHEN | From 426 BCE to the present |
| WHERE | All over the earth |



**ANSWER:** The histogram displays the distribution of earthquake magnitudes on the Richter scale. The height of the tallest bar says that there were about 250 earthquakes with magnitudes between 7.0 and 7.5. We can see that earthquakes typically have magnitudes around 7. Most are between 5.5 and 8.5, but one is less than 4 and a few are 9 or bigger. Relative to the other tsunami-causing earthquakes, the Sumatra and Japan events were extraordinarily powerful.

## EXAMPLE 2.4

## How Much Do Americans Work?

The Bureau of Labor Statistics (BLS) collects data on many aspects of the U.S. economy. One of the surveys it conducts, the American Time Use Survey (ATUS), asks roughly 11,000 people a year a variety of questions about how they spend their time. For those who are employed, it asks how many hours a week they work. Here is a histogram of the 2270 responses in 2014.

**QUESTION:** What does the histogram say about how many hours U.S. workers typically work?



**ANSWER:** It looks like the vast majority of people (more than 1200 in this study) work right around 40 hours a week. There are some who work less, and a very few who work more.

# Stem-and-Leaf Displays

Histograms provide an easy-to-understand summary of the distribution of a quantitative variable, but they don't show the data values themselves. For example, here's a histogram of the pulse rates of 24 women at a health clinic:

**Figure 2.8**

A histogram of the pulse rates of 24 women at a health clinic.



Here's a stem-and-leaf display of the same data:

```
5 | 6
6 | 0 4 4 4
6 | 8 8 8 8
7 | 2 2 2 2
7 | 6 6 6 6
8 | 0 0 0 0 4 4
8 | 8
Pulse Rate
(5|6 means 56 beats/min)
```
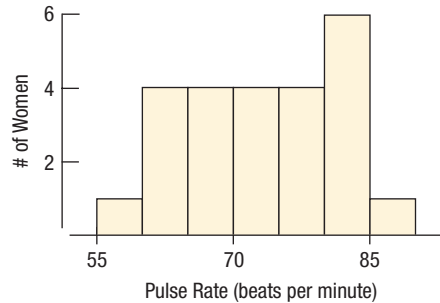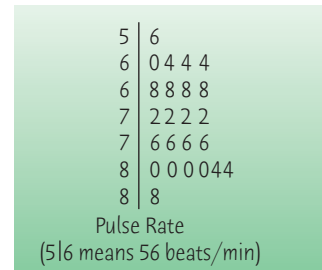
**STEM-AND-LEAF OR STEMPLOT?**

The stem-and-leaf display was devised by John W. Tukey, one of the greatest statisticians of the 20th century. It is called a "stemplot" in some texts and computer programs.

A **stem-and-leaf display** is like a histogram, but it shows the individual values. It's also easier to make by hand. Turn the stem-and-leaf on its side (or turn your head to the right) and squint at it. It should look roughly like the histogram of the same data. Does it?[5]

The first line of the display, which says 5│6, stands for a pulse of 56 beats per minute (bpm). We've taken the tens place of the number and made that the "stem." Then we sliced off the ones place and made it a "leaf." The next line down is 6│0444, which shows one pulse rate of 60 and three of 64 bpm.

Stem-and-leaf displays are especially useful when you make them by hand for batches of fewer than a few hundred data values. They are a quick way to display—and even to record—numbers. Because the leaves show the individual values, we can sometimes see even more in the data than the distribution's shape. Take another look at all the leaves of the pulse data. See anything unusual? At a glance you can see that they are all even. With a bit more thought you can see that they are all multiples of 4—something you couldn't possibly see from a histogram. How do you think the nurse took these pulses? Counting beats for a full minute, or counting for only 15 seconds and multiplying by 4?

# Dotplots

A **dotplot** places a dot along an axis for each case in the data. It's like a stem-and-leaf display, but with dots instead of digits for all the leaves. Dotplots are a great way to display a small dataset. Figure 2.9 shows a dotplot of the time (in seconds) that the winning horse took to win the Kentucky Derby in each race between the first Derby in 1875 and the 2018 Derby. (Data in **Kentucky Derby 2018**)
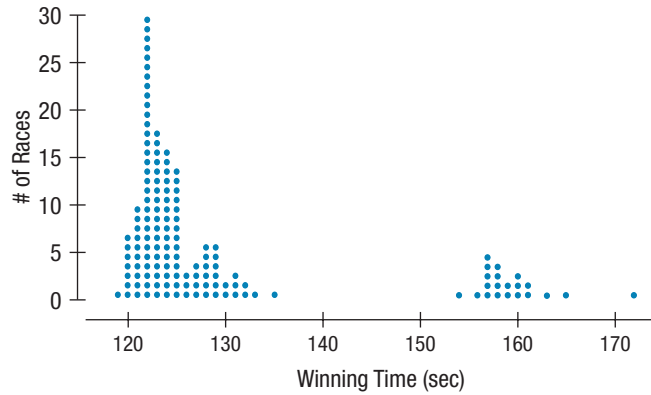
---

[5]You could make the stem-and-leaf display with the higher values on the top. Putting the lower values at the top matches the histogram; putting the higher values at the top matches the way a vertical axis works in other displays such as dotplots (as we'll see presently).

Dotplots display basic facts about the distribution. We can find the slowest and fastest races by finding the times for the topmost and bottommost dots. It's clear that there are two clusters of points, one just below 160 seconds and the other at about 122 seconds. Something strange happened to the Derby times. Once we know to look for it, we can find out that in 1896 the distance of the Derby race was changed from 1.5 miles to the current 1.25 miles. That explains the two clusters of winning times.

**Figure 2.9**

A dotplot of Kentucky Derby winning times plots each race as its own dot. We can see two distinct groups corresponding to the two different race distances.

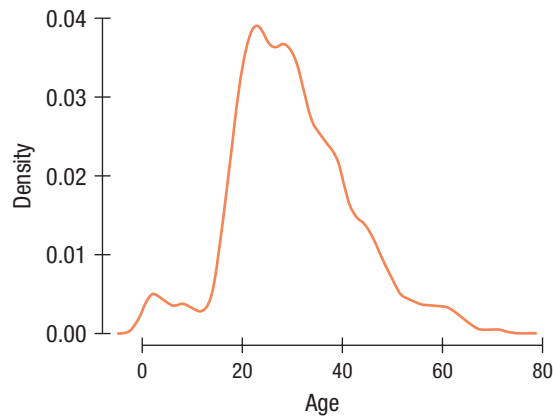| WHO | Runnings of the Kentucky Derby |
|---|---|
| WHAT | Winning time |
| WHEN | 1875–2018 |
| WHERE | Churchill Downs |



## *Density Plots

The size of the bins in a histogram can influence its look and our interpretation of the distribution. There is no correct bin size, although recommendations to use between 5 and 20 bins are common. **Density plots** smooth the bins in a histogram to reduce the effect of this choice. How much the bin heights are smoothed is still a choice that affects the shape, but the change in shape is less severe than in a histogram. Here's a density plot of the *Ages* of those on the *Titanic*. Compare it to Figure 2.7.

**Figure 2.10**

A density plot of the *Ages* of those aboard the *Titanic*. We can see, as we did in the histogram, that the most populous age is near 20 and that there are more infants and toddlers than preteens. The density plot does not provide hard cut-offs to the bins, but smooths the distribution over the bins.



Every histogram, stem-and-leaf display, and dotplot tells a story, but you need to develop a vocabulary to help you explain it. Start by talking about three things: its *shape*, *center*, and *spread*.

### Think Before You Draw

Before making a pie chart or a bar chart, you should check that you have categorical data. Before making a stem-and-leaf display, a histogram, or a dotplot, you should make sure you are working with quantitative data. Although a bar chart and a histogram may look similar, they're not the same display. You can't display categorical data in a histogram nor quantitative data in a bar chart.

## 2.3 Shape

We summarize the **shape** of a distribution in terms of three attributes: how many *modes* it has, whether it is *symmetric* or *skewed*, and whether it has any extraordinary cases or *outliers*.

1.  *Does the histogram have a single, central hump or several separated humps?* These humps are called **modes**.[6] The mode is sometimes defined as the single value that appears most often. That definition is fine for categorical variables because all we need to do is count the number of cases for each category. For quantitative variables, the mode is more ambiguous. It makes more sense to use the term "mode" to refer to the peak of the histogram rather than as a single summary value. The important feature of the Kentucky Derby races is that there are two distinct modes, representing the two different versions of the race and warning us to consider those two versions separately. The earthquake magnitudes of Example 2.3 have a single mode at just about 7.

    A histogram with one peak, such as the ages (Figure 2.7), is dubbed **unimodal**; histograms with two peaks are **bimodal**, and those with three or more are called **multimodal**.[7]
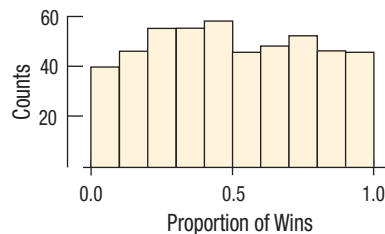
---

*PIE À LA MODE?*

You've heard of pie à la mode. Is there a connection between pie and the mode of a distribution? Actually, there is! The mode of a distribution is a *popular* value near which a lot of the data values gather. And "à la mode" means "in style"—*not* "with ice cream." That just happened to be a *popular* way to have pie in Paris around 1900.

---

A histogram that doesn't appear to have any mode and in which all the bars are approximately the same height is called **uniform**.
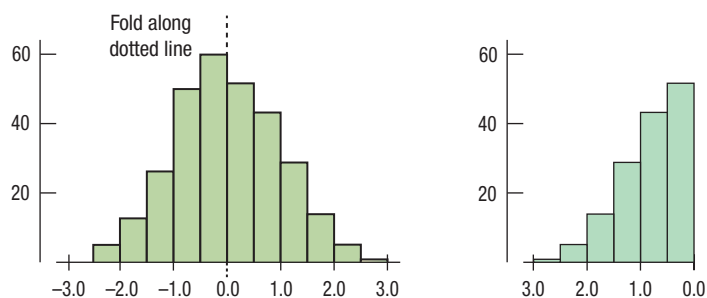
**Figure 2.11**

In this histogram, the bars are all about the same height. The histogram doesn't appear to have a mode and is called uniform.



2.  *Is the histogram* **symmetric**? Can you fold it along a vertical line through the middle and have the edges match pretty closely, or are more of the values on one side?

**Figure 2.12**

A symmetric histogram can fold in the middle so that the two sides almost match.



---

[6]Well, technically, it's the value on the horizontal axis of the histogram that is the mode, but anyone asked to point to the mode would point to the hump.

[7]Apparently, statisticians don't like to count past two.