AP STATISTICS - SUMMER ASSIGNMENT

For the AP® Exam

The Practice of Statistics

FIFTH EDITION

PLEASE DO NOT MARK
ON COPY AND TURN
IN AT BEGINNING
OF SCHOOL YEAR

MR. SCHAUER

To the Student

Statistical Thinking and You

The purpose of this book is to give you a working knowledge of the big ideas of statistics and of the methods used in solving statistical problems. Because data always come from a real-world context, doing statistics means more than just manipulating data. The Practice of Statistics (TPS), Fifth Edition, is full of data. Each set of data has some brief background to help you understand what the data say. We deliberately chose contexts and data sets in the examples and exercises to pique your interest.

TPS 5e is designed to be easy to read and easy to use. This book is written by current high school AP® Statistics teachers, for high school students. We aimed for clear, concise explanations and a conversational approach that would encourage you to read the book. We also tried to enhance both the visual appeal and the book's clear organization in the layout of the pages.

Be sure to take advantage of all that TPS 5e has to offer. You can learn a lot by reading the text, but you will develop deeper understanding by doing Activities and Data Explorations and answering the Check Your Understanding questions along the way. The walkthrough guide on pages xiv—xx gives you an inside look at the important features of the text.

You learn statistics best by doing statistical problems. This book offers many different types

of problems for you to tackle.

- Section Exercises include paired odd- and even-numbered problems that test the same skill or concept from that section. There are also some multiple-choice questions to help prepare you for the AP[®] exam. Recycle and Review exercises at the end of each exercise set involve material you studied in previous sections.
- Chapter Review Exercises consist of free-response questions aligned to specific learning objectives from the chapter. Go through the list of learning objectives summarized in the Chapter Review and be sure you can say "I can do that" to each item. Then prove it by solving some problems.
- The AP® Statistics Practice Test at the end of each chapter will help you prepare for in-class exams. Each test has 10 to 12 multiple-choice questions and three free-response problems, very much in the style of the AP® exam.
- Finally, the Cumulative AP® Practice Tests after Chapters 4, 7, 10, and 12 provide challenging, cumulative multiple-choice and free-response questions like ones you might find on a midterm, final, or the AP® Statistics exam.

The main ideas of statistics, like the main ideas of any important subject, took a long time to discover and take some time to master. The basic principle of learning them is to be persistent. Once you put it all together, statistics will help you make informed decisions based on data in your daily life.

TPS and AP® Statistics

£

1...

:8

u

d

;h

a-

be

3.

1e

ıd

วท

es

S-

ach

ſе

:e-

ride

/ou

ie to

tent.

ta in

The Practice of Statistics (TPS) was the first book written specifically for the Advanced Placement (AP®) Statistics course. Like the previous four editions, TPS 5e is organized to closely follow the AP® Statistics Course Description. Every item on the College Board's "Topic Outline" is covered thoroughly in the text. Look inside the front cover for a detailed alignment guide. The few topics in the book that go beyond the AP® syllabus are marked with an asterisk (*).

Most importantly, TPS 5e is designed to prepare you for the AP® Statistics exam. The entire author team has been involved in the AP® Statistics program since its early days. We have more than 80 years' combined experience teaching introductory statistics and more than 30 years' combined experience grading the AP® exam! Two of us (Starnes and Tabor) have served as Question Leaders for several years, helping to write scoring rubrics for free-response questions. Including our Content Advisory Board and Supplements Team (page vii), we have two former Test Development Committee members and 11 AP® exam Readers.

TPS 5e will help you get ready for the AP® Statistics exam throughout the course by:

- Using terms, notation, formulas, and tables consistent with those found on the AP® exam. Key terms are shown in bold in the text, and they are defined in the Glossary. Key terms also are cross-referenced in the Index. See page F-1 to find "Formulas for the AP® Statistics Exam" as well as Tables A, B, and C in the back of the book for reference.
- Following accepted conventions from AP® exam rubrics when presenting model solutions. Over the years, the scoring guidelines for free-response questions have become fairly consistent. We kept these guidelines in mind when writing the solutions that appear throughout TPS 5e. For example, the four-step State-Plan-Do-Conclude process that we use to complete inference problems in Chapters 8 through 12 closely matches the four-point AP® scoring rubrics.
- Including AP® Exam Tips in the margin where appropriate. We place exam tips in the margins and in some Technology Corners as "on-the-spot" reminders of common mistakes and how to avoid them. These tips are collected and summarized in Appendix A.
- Providing hundreds of AP®-style exercises throughout the book. We even added a new kind of problem just prior to each Chapter Review, called a FRAPPY (Free Response AP® Problem, Yay!). Each FRAPPY gives you the chance to solve an 'AP®-style free-response problem based on the material in the chapter. After you finish, you can view and critique two example solutions from the book's Web site (www.whfreeman.com/tps5e). Then you can score your own response using a rubric provided by your teacher.

Turn the page for a tour of the text. See how to use the book to realize success in the course and on the AP® exam.

READ THE TEXT and use the book's features to help you grasp the big ideas.

Read the LEARNING
OBJECTIVES at the
beginning of each section.
Focus on mastering these
skills and concepts as you
work through the chapter.

Scan the margins for the purple notes, which represent the "voice of the teacher" giving helpful hints for being successful in the course.

Look for the boxes with the blue bands. Some explain how to make graphs or set up calculations while others recap important concepts.

Often, using the regression line

to make a prediction for x = 0 k

an extrapolation. That's why the y

rcept isn't always statistically

3.1 Scatterplots and Correlation

WHAT YOU WILL LEARN

By the end of the section, you should be able to:

- Identify explanatory and response variables in situations where one variable helps to explain or influences the other.
- Make a scatterplot to display the relationship between two quantitative variables.
- Describe the direction, form, and strength of a relationship displayed in a scatterplot and identify outliers in a scatterplot.
- Interpret the correlation.
- Understand the basic properties of correlation, including how the correlation is influenced by outliers.
- Use technology to calculate correlation.
- Explain why association does not imply causation.

DEFINITION: Extrapolation

Extrapolation is the use of a regression line for prediction far outside the interval of values of the explanatory variable x used to obtain the line. Such predictions are often not accurate.

Few relationships are linear for all values of the explanatory variable. Don't make predictions using values of x that are much larger or much smaller than those that actually appear in your data.



Take note of the green DEFINITION boxes that explain important vocabulary. Flip back to them to review key terms and their definitions.

Watch for CAUTION ICONS. They alert you to common mistakes that students make.

HOW TO MAKE A SCATTERPLOT

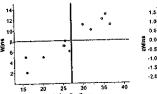
- Decide which variable should go on each axis.
- 2. Label and scale your axes.
- 3. Plot individual data values.

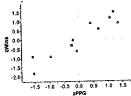
Make connections and deepen your understanding by reflecting on the questions asked in THINK ABOUT IT passages.

THINK ABOUT IT

What does correlation measure? The Fathom screen shots below provide more detail. At the left is a scatterplot of the SEC football data with two lines added—a vertical line at the group's mean points per game and a horizontal line at the mean number of wins of the group. Most of the points fall in the upper-right or lower-left "quadrants" of the graph. That is, teams with above-average points per game tend to have above-average numbers of wins, and teams with below-average points per game tend to have numbers of wins that are below average. This confirms the positive association between the variables.

Below on the right is a scatterplot of the standardized scores. To get this graph, we transformed both the x- and the y-values by subtracting their mean and dividing by their standard deviation. As we saw in Chapter 2, standardizing a data set converts the mean to 0 and the standard deviation to 1. That's why the vertical and horizontal lines in the right-hand graph are both at 0.





Notice that all the products of the standardized values will be positive—not surprising, considering the strong positive association between the variables. What if there was a negative association between two variables? Most of the points would be in the upper-left and lower-right "quadrants" and their z-score products would be negative, resulting in a negative correlation.

Read the AP® EXAM TIPS. They give advice on how to be successful on the AP® exam. AP® EXAM TIP The formula sheet for the AP® exam uses different notation for these equations: $b_1 = r \frac{S_r}{S_s}$ and $b_0 = \overline{y} - b_1 \overline{x}$. That's because the least-squares line is written as $\beta = b_0 + b_1 x$. We prefer our simpler versions without the subscripts!



LEARN STATISTICS BY DOING STATISTICS



ACTIVITY Reaching for Chips

MATERIALS:

200 colored chips, including 100 of the same color; large bag or other container



Before class, your teacher will prepare a population o having the same color (say, red). The parameter is this in the population: $\rho = 0.50$. In this Activity, y variability by taking repeated random samples of size

 After your teacher has mixed the chips thoroughly should take a sample of 20 chips and note the sample When finished, the student should return all the chip and pass the bag to the next student.

Note: If your class has fewer than 25 students, have s samples.

- 2. Each student should record the \hat{p} -value in a chart value on a class dotplot. Label the graph scale from f spaced 0.05 units apart.
- 3. Describe what you see: shape, center, spread, and usual features.

Every chapter begins with a hands-on ACTIVITY that introduces the content of the chapter. Many of these activities involve collecting data and drawing conclusions from the data. In other activities, you'll use dynamic applets to explore statistical concepts.

DATA EXPLORATIONS ask you to play the role of data detective. Your goal is to answer a puzzling, real-world question by examining data graphically and numerically.

ACTIVITY I'm a Great Free-Throw Shooter!

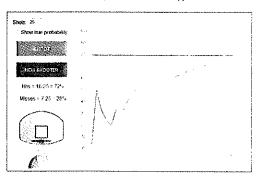
MATERIALS:

Computer with Internet access and projection capability



A basketball player claims to make 80% of the free throws that he attempts. We think he might be exaggerating. To test this claim, we'll ask him to shoot some free throws—virtually—using The Reasoning of a Statistical Test applet at the book's Web site.

1. Go to www.whfreeman.com/tps5e and launch the applet.



- 2. Set the applet to take 25 shots. Click "Shoot." How many of the 25 shots did the player make? Do you have enough data to decide whether the player's claim is valid?
- 3. Click "Shoot" again for 25 more shots. Keep doing this until you are convinced either that the player makes less than 80% of his shots or that the player's claim is true. How large a sample of shots did you need to make your decision?
- 4. Click "Show true probability" to reveal the truth. Was your conclusion correct?
- 5. If time permits, choose a new shooter and repeat Steps 2 through 4. Is it easier to tell that the player is exaggerating when his actual proportion of free throws made is closer to 0.8 or farther from 0.8?

THUL WELLOWAYOW The SAT essay: Is longer better?

Following the debut of the new SAT Writing test in March 2005, Dr. Les Perelman from the Massachusetts Institute of Technology stirred controversy by reporting, "It appeared to me that regardless of what a student wrote, the longer the essay, the ligher the score." He went on to say, "I have never found a quantifiable predictor in 25 years of grading that was anywhere as strong as this one. If you just graded them based on length without ever reading them, you'd be right over 90 percent of the time." The table below shows the data that Dr. Perelman used to draw his conclusions."

| 小数 變 | dema | L. | offi of s | | 6 200 78 | Sar 8 🗪 | क्रकृष्टि औ | SAT ESS | 100 | Arrest Marie | |
|-------------|------|-----|-----------|-----|-----------------|---------|-------------|---------|-----|-----------------|-----|
| Words: | 460 | 422 | 402 | 365 | 357 | 278 | 236 | 201 | 168 | 156 | 133 |
| Score: | 6 | 6 | 5 | 5 | 6 | 5 | 4 | 4 | 4 | 3 | 2 |
| Words: | 114 | 108 | 100 | 403 | 401 | 388 | 320 | 258 | 236 | 189 | 128 |
| Score: | 2 | 1 | 1 | 5 | 6 | 6 | 5 | 4 | 4 | 3 | 2 |
| Words: | 67 | 697 | 387 | 355 | 337 | 325 | 272 | 150 | 135 | | |
| Score: | 1 | 6 | 6 | 5 | 5 | 4 | 4 | 2 | 3 | | |

Does this mean that if students write a lot, they are guaranteed high scores? Carry out your own analysis of the data. How would you respond to each of Dr. Perelman's claims?

CHECK YOUR UNDERSTANDING questions appear throughout the section. They help you to clarify definitions, concepts, and procedures. Be sure to check your answers in the back of the book.



CHECK YOUR UNDERSTANDING

Identify the explanatory and response variables in each setting

- 1. How does drinking beer affect the level of alcohol in people's blood? The legal limit for driving in all states is 0.08%. In a study, adult volunteers drank different numbers of cans of beer. Thirty minutes later, a police officer measured their blood alcohol levels.
- 2. The National Student Loan Survey provides data on the amount of debt for recent college graduates, their current income, and how stressed they feel about college debt. A sociologist looks at the data with the goal of using amount of debt and income to explain the stress caused by college debt.

EXAMPLES: Model statistical problems and how to solve them

CHAPTER 3 DESCRIBING RELATIONSHIPS

You will often see explanatory variables called independent variables and response variables called dependent variables. Recause the words 'independent" and "dependent" have other meanings in statistics, we won't use them here.

It is easiest to identify explanatory and response variables when we actually specify values of one variable to see how it affects another variable. For instance, to study the effect of alcohol on body temperature, researchers gave several different amounts of alcohol to mice. Then they measured the change in each

mouse's body temperature 15 minutes later. In this case, amount of alcohol is the explanatory variable, and change in body temperature is the response variable. When we don't specify the values of either variable but just observe both variables, there may or may not be explanatory and response variables. Whether there are depends on how you plan to use the data.

Read through each EXAMPLE, and then try out the concept yourself by working the FOR PRACTICE exercise in the



Linking SAT Math and Critical Reading Scores

Explanatory or response?

Julie asks, "Can I predict a state's mean SAT Math score if I know its mean SAT Critical Reading score?" Jim wants to know how the mean SAT Math and Critical Reading scores this year in the 50 states are related to each other.

PROBLEM: For each student, identify the explanatory variable and the response variable if possible. SOLUTION: Julie is treating the mean SAT Critical Reading score as the explanatory variable and the mean SAT Math score as the response variable. Jim is simply interested in exploring the relationship between the two variables. For him, there is no clear explanatory or response variable.

For Practice Try Exercise

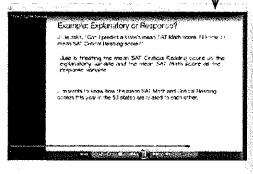
Section Exercises. Need extra help? Examples

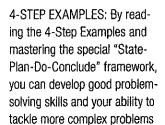
and exercises marked with the PLAY ICON (are supported by short video clips prepared by experienced AP® teachers. The video guides you through each step in the example and solution and gives you extra help when you need it.

The red number box next to the exercise directs you back to the page in the section where the model example appears.



Coral reefs How sensitive to changes in water temperature are coral reefs? To find out, measure the growth of corals in aquariums where the water temperature is controlled at different levels. Growth is measured by weighing the coral before and after the experiment. What are the explanatory and response variables? Are they categorical or quantitative?





like those on the AP® exam.





Gesell Scores

Putting it all together

Does the age at which a child begins to talk predict a later score on a test of mental ability? A study of the development of young children recorded the age in months at which each of 21 children spoke their first word and their Gesell Adaptive Score, the result of an aptitude test taken much later. The data appear in the table below, along with a scatterplot, residual plot, and computer output. Should we use a linear model to predict a child's Gesell score from his or her age at first word? If so, how accurate will our predictions be?

| | NOTE: | Age | months) | at first w | ord and Ge | sell score | | 经基础的 |
|-------|-------|-------|---------|------------|------------|------------|-----|-------|
| CHILD | AGE | SCORE | CHILD | AGE | SCORE | CHILD | AGE | SCORE |
| 1 | 15 | 95 | 8 | 11 | 100 | 15 | 11 | 102 |
| 2 | 26 | 71 | 9 | 8 | 104 | 16 | 10 | 100 |
| 3 | 10 | 83 | 10 | 20 | 94 | 17 | 12 | 105 |
| 4 | 9 | 91 | 11 | 7 | 113 | 18 | 42 | 57 |
| 5 | 15 | 102 | 12 | 9 | 96 | 19 | 17 | 121 |
| 6 | 20 | 87 | 13 | 19 | 83 | 20 | 17 | 86 |
| 7 | 18 | 93 | 14 | 11 | 84 | 21 | 10 | 100 |

EXERCISES: Practice makes perfect!



Section 3.2 Summary

- A regression line is a straight line that describes how a response variable y changes
 as an explanatory variable x changes. You can use a regression line to predict the
 value of y for any value of x by substituting this x into the equation of the line.
- The slope b of a regression line ŷ = a + bx is the rate at which the predicted
 response ŷ changes along the line as the explanatory variable x changes. Specifically, b is the predicted change in y when x increases by 1 unit.
- The y intercept a of a regression line ŷ = a + bx is the predicted response ŷ when the explanatory variable x equals 0. This prediction is of no statistical use unless x can actually take values near 0.

Start by reading the SECTION SUMMARY to be sure that you understand the key concepts.

Practice! Work the EXERCISES assigned by your teacher. Compare your answers to those in the Solutions Appendix at the back of the book. Short solutions to the exercises numbered in red are found in the appendix.

Most of the exercises are paired, meaning that odd- and even-numbered problems test the same skill or concept. If you answer an assigned problem incorrectly, try to figure out your mistake. Then see if you can solve the paired exercise.

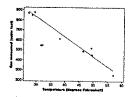
Look for icons that appear next to selected problems. They will guide you to

- an Example that models the problem.
- videos that provide stepby-step instructions for solving the problem.
- earlier sections on which the problem draws (here, Section 2.2).
- examples with the
 4-Step State-Plan-Do-Conclude way of solving problems.

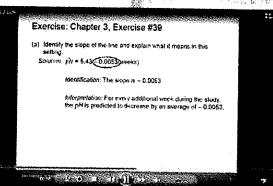
Section 3.2 Exercises

- 35. What's my line? You use the same bar of soap to shower each morning. The bar weighs 80 grams when it is new. Its weight goes down by 6 grams per day on average. What is the equation of the regression line for predicting weight from days of use?
- 36. What's my line? An eccentric professor believes that a child with IQ 100 should have a reading test score of 50 and predicts that reading score should increase by 1 point for every additional point of IQ. What is the equation of the professor's regression line for predicting reading score from IQ?
- 37. Gas mileage We expect a car's highway gas mileage to be related to its city gas mileage. Data for all 1198 vehicles in the government's recent Fuel Economy Guide give the regression line: predicted highway mpg = 4.62 + 1.109 (city mpg).
- (a) What's the slope of this line? Interpret this value in context
- (b) What's the y intercept? Explain why the value of the intercept is not statistically meaningful.
- (c) Find the predicted highway mileage for a car that gets 16 miles per gallon in the city.
- 38. IQ and reading scores Data on the IQ test scores and reading test scores for a group of fifth-grade children give the following regression line: predicted reading score = −33.4 + 0.882(IQ scog. ...
- (a) What's the slope of this line? Interpret this context.
- (b) What's the y intercept? Explain why the intercept is not statistically meaningful.
- (c) Find the predicted reading score for a cf IQ score of 90.
 - 9. Acid rain Researchers studying acid rain the acidity of precipitation in a Colorade area for 150 consecutive weeks. Acidity in the property of the consecutive weeks. Acidity in the property of the proper

in Joan's midwestern home. The figure below shows the original scatterplot with the least-squares line added. The equation of the least-squares line is $\hat{y} = 1425 - 19.87x$.



- (a) Identify the slope of the line and explain what it means in this setting.
- (b) Identify the y intercept of the line. Explain why it's risky to use this value as a prediction.
- (c) Use the regression line to predict the amount of natural gas Joan will use in a month with an average temperature of 30°F.
- Acid rain Refer to Exercise 39. Would it be appropriate to use the regression line to predict pH after 1000 months? Justify your answer.



79. In my Chevrolet (2.2) The Chevrolet Malibu with

a four-cylinder engine has a combined gas mileage of

25 mpg. What percent of all vehicles have worse gas
mileage than the Malibu?

67. Beavers and beetles Do beavers benefit beetles?
4. Researchers laid out 23 circular plots, each 4 meters in diameter, in an area where beavers were cutting down cottonwood trees. In each plot, they counted the number of stumps from trees cut by beavers and the number of clusters of beetle larvae. Ecologists think that the new sprouts from stumps are more tender than other cottonwood growth, so that beetles prefer them.



REVIEW and PRACTICE for quizzes and tests



Chapter Review

Section 3.1: Scatterplots and Correlation

In this section, you learned how to explore the relationship between two quantitative variables. As with distributions of a single variable, the first step is always to make a graph. A scatterplot is the appropriate type of graph to investigate associations between two quantitative variables. To describe a scatterplot, be sure to discuss four characteristics direction, form, strength, and outliers. The direction of an association might be positive, negative, or neither. The form of an association can be linear or nonlinear. An association is strong if it closely follows a specific form. Finally, outliers are any points that clearly fall outside the pattern of the rest of the data.

The correlation r is a numerical summary that describes the direction and strength of a linear association. When r > 0, the association is positive, and when r < 0, the association is negative. The correlation will always take values between -1 and 1, with r = -1 and r = 1 indicing a perfectly linear relationship. Strong linear associations have correlations near 1 or -1, while weak linear relationships have correlations near 0. However, it is

Use the WHAT DID YOU LEARN? table to guide you to model examples and exercises to verify your mastery of

each LEARNING OBJECTIVE.

possible to determine the form of an association from only the correlation. Strong nonlinear relationships can have a correlation close to 0 or a correlation close to 0, depending on the association. You also learned that outliers can greatly affect the value of the correlation and that correlation does not imply causation. That is, we can't assume that changes in one variable cause changes in the other variable, just because they have a correlation close to 1 or -1.

Section 3.2: Least-Squares Regression

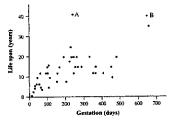
In this section, you learned how to use least-squares regression lines as models for relationships between variables that have a linear association. It is important to understand the difference between the actual data and the model used to describe the data. For example, when you are interpreting the slope of a least-squares regression Review the CHAPTER SUMMARY to be sure that you understand the key concepts in each section.

| Learning Objective | Section | Related Example on Page(s) | Relevant Chapter Review Exercise(s |
|---|---------|---|---------------------------------------|
| dentify explanatory and response variables in situations where one variable heips to explain or influences the other. | 3.1 | ` 144 | R3.4 |
| Make a scatterplot to display the relationship between two quantitative variables. | 3.1 | 145, 148 | R3.4 |
| Describe the direction, form, and strength of a relationship displayed in a scatterplot and recognize outliers in a scatterplot. | 3.1 | 147, 148 | R3.1 |
| Interpret the correlation. | 3.1 | 152 | R3.3, R3.4 |
| Understand the basic properties of correlation, including how the correlation is influenced by outliers. | 3.1 | 152, 156, 157 | 83.1, 83.2 |
| Use technology to calculate correlation. | 3.1 | Activity on 152, 171 | R3.4 |
| Explain why association does not imply causation. | 3.1 | Discussion on 156, 190 | R3.6 |
| Interpret the slope and y intercept of a least-squares regression line. | 3.2 | 166 | R3.2, R3.4 |
| Use the least-squares regression line to predict y for a given x. Explain the dangers of extrapolation. | 3.2 | 167, Discussion on 168 (for extrapolation) | R3.2, R3.4, R3.5 |
| Calculate and interpret residuals. | 3.2 | 169 | R3.3, R3.4 |
| Explain the concept of least squares. | 3.2 | Discussion on 169 | R3.5 |
| Determine the equation of a least-squares regression line using technology or computer output. | 3.2 | Technology Corner on 171, 181 | R3.3, R3.4 |
| Construct and Interpret residual plots to assess whether a linear model is appropriate. | 3.2 | Discussion on 175, 180 | R3.3, R3.4 |
| Interpret the standard deviation of the residuals and r ² and use these values to assess how well the least-squares regression line models the relationship between two variables. | 3.2 | 180 | R3.3, R3.5 |
| | | Discussion on 188 | R3.1 |

Chapter 3 Chapter Review Exercises

These exercises are designed to help you review the important ideas and methods of the chapter.

R3.1 Born to be old? Is there a relationship between the gestational period (time from conception to birth) of an animal and its average life span? The figure shows a scutterplot of the gestational period and average life span for 43 species of animals.³⁰



(a) Describe the association shown in the scatterplot.

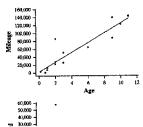
R3.3 Stats teachers' cars A random sample of AP[®] Statistics teachers was asked to report the age (in years) and mileage of their primary vehicles. A scatterplot of the data, a least-squares regression printout, and a residual plot are provided below.

 Predictor
 Coef
 SE Coef
 T
 P

 Constant
 3704
 8259
 0.45
 0.66

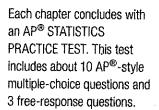
 Age
 12188
 1492
 8.17
 0.00

S = 20870.5 R-Sq = 83.7% R-Sq(adj) = 82.4%



Tackle the CHAPTER REVIEW EXERCISES for practice in solving problems that test concepts from throughout the chapter.

and the AP® Exam

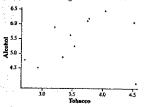


Chapter 3 AP® Statistics Practice Test

Section I: Multiple Chains Select the host answer for each question

- T3.1 A school guidance counselor examines the number of extracurricular activities that students do and their grade point average. The guidance counselor says, "The evidence indicates that the correlation between the number of extracurricular activities a student participates in and his or her grade point average is close to zero." A correct interpretation of this statement would be that
- (a) active students tend to be students with poor grades, and vice versa.
- (b) students with good grades tend to be students who are not involved in many extracurricular activities, and vice versa.
- (c) students involved in many extracumoular activities are just as likely to get good grades as bad grades; the same is true for students involved in few extracumoular activities.
- (d) there is no linear relationship between number of activ-

alcoholic beverages for each of 11 regions in Great Britain was recorded. A scatterplot of spending on alcohol versus spending on tobacco is shown below. Which of the following statements is true?



- (a) The observation (4.5, 6.0) is an outlier.
- There is clear evidence of a negative association beending on alcohol and tobacco.

 action of the least-squares line for this plot
 be approximately \$\tilde{y} = 10 - 2x\$,
 celation for these data is \$r \neq 0.99,
 struction in the lower-right corner of the plot is
 all for the least-squares line.

Cumulative AP® Practice Test 1

Section t: Multiple Choice Choose the best answer for Questions AP1.1 to AP1.14.

- AP1.1 You look at real estate ads for houses in Sarasota. Florida. Many houses range from \$200,000 to \$400,000 in price. The few houses on the water, however, have prices up to \$15 million. Which of the following statements best describes the distribution of home prices in Sarasota?
 - (a) The distribution is most likely skewed to the left, and the mean is greater than the median.

AP1.4 For a certain experiment, the available experimental units are eight rats, of which four are female (F1, F2, F3, F4) and four are male (M1, M2, M3, M4). There are to be four treatment groups, A, B, C, and D. If a randomized block design is used, with the experimental units blocked by gender, which of the following assignments of treatments is impossible?

Section II: Free Response Show all your work. Indicate clearly the methods you use, because you will be graded on the correctness of your methods as well as on the accuracy and completeness of your methods as well as on the accuracy and completeness of your results and explanations.

AP1.15 The manufacturer of exercise machines for fitness centers has designed two new elliptical machines that are meant to increase cardiovascular fitness. The two machines are being tested on 30 volunteers at a fitness center near the company's headquarters. The volunteers are randomly assigned to one of the machines and use it daily for two months. A measure of cardiovascular fitness is administered at the start of the experiment and

the two machines. Note that higher scores indicate larger gains in fitness.

| Machine A | | Machine E |
|-----------|--------|-----------|
| | 0 | 2 |
| 54 | 1 | 0 |
| 876320 | 2 | 159 |
| 97411 | 3 | 2489 |
| | ا بر ا | 257 |

Four CUMULATIVE AP® TESTS simulate the real exam. They are placed after Chapters 4, 7, 10, and 12. The tests expand in length and content coverage from the first through the fourth.

Learn how to answer free-response

the FRAPPY! THE FREE RESPONSE AP® PROBLEM, YAY! that comes just before the Chapter Review in every chapter.

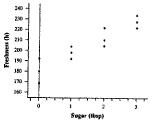
questions successfully by working

FRAPPY! Free Response AP® Problem, Yay!

The following problem is modeled after actual AP[®] Statistics exam free response questions. Your task is to generate a complete, concise reasonse in 15 minutes

Directions: Show all your work. Indicate clearly the methods you use, because you will be scored on the correctness of your methods as well as on the accuracy and completeness of your results and explanations.

and observed how many hours each flower continued to look fresh. A scatterplot of the data is shown below.



- (a) Briefly describe the association shown in the scatterplot.
- (b) The equation of the least-squares regression line for these data is ŷ = 180.8 + 15.8x. Interpret the slope of the line in the context of the study.

Two statistics students went to a flower shop and randomly selected 12 carnations. When they got home, the students prepared 12 identical vases with exactly the same amount of water in each vase. They put one tablespoon of sugar in 3 vases, two tablespoons of sugar in 3 vases, and three tablespoons of sugar in 3 vases. In the remaining 3 vases, they put no sugar. After the vases were prepared, the students randomly assigned 1 carnation to each vase

- (c) Calculate and interpret the residual for the flower that had 2 tablespoons of sugar and looked fresh for 204 hours.
- (d) Suppose that another group of students conducted a similar experiment using 12 flowers, but included different varieties in addition to camations. Would you expect the value of r' for the second group's data to be greater than, less than, or about the same as the value of r' for the first group's data? Explain.

After you finish, you can view two example solutions on the book's Web site (www.wifreeman.com/tps5e). Determine whether you think each solution is "complete," "substantial," "developing," or "minimal." if the solution is not complete, what improvements would you suggest to the student who wrote It? Finally, your teacher will provide you with a scoring rubric. Score your response and note what, if anything, you would do differently to improve your own score.

Use TECHNOLOGY to discover and analyze



Use technology as a tool for discovery and analysis. TECHNOLOGY CORNERS give step-by-step instructions for using the TI-83/84 and TI-89 calculator. Instructions for the TI-Nspire are in an end-of-book appendix. HP Prime instructions are on the book's Web site and in the e-Book.

TECHNOLOGY CORNER

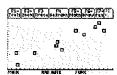
Making scatterplots with technology is much easier than constructing them by hand. We'll use the SEC football data from page 146 to show how to construct a scatterplot on a TI-83/84 or TI-89.

- Enter the data values into your lists. Put the points per game in L1/list1 and the number of wins in L2/list2.
- Define a scatterplot in the statistics plot menu (press F2 on the TI-89). Specify the settings shown below.



 $\label{thm:comData} Use Zoom Stat (Zoom Data on the T1-89) to obtain a graph. The calculator will set the window dimensions automatically the compact of t$ by looking at the values in L1/list1 and L2/list2.





Notice that there are no scales on the axes and that the axes are not labeled. If you copy a scatterplot from your calculator onto your paper, make sure that you scale and label the axes

> AP® EXAM TIP. If you are asked to make a scatterplot on a free-response question, be sure to label and scale both axes. Don't just copy an unlabeled calculator graph directly onto your paper

You can access video instructions for the Technology Corners through the e-Book or on the book's Web site.



Find the Technology Corners easily by consulting the summary table at the end of each section or the complete table inside the back cover of the book.



9. Residual plots on the calculator

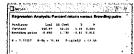
8. Least-squares regression lines on the calculator

page 171 page 175





59. Merlins breeding Exercise 13 (page 160) gives data on the number of breeding pairs of merlins in an isolated area in each of seven years and the percent of males who returned the next year. The data show that the percent returning is lower after successful breeding seasons and that the relationship is roughly linear. The figure below shows Minitab regression output for these data.



- (a) What is the equation of the least-squares regression line for predicting the percent of males that return from the number of breeding pairs? Use the equation to predict the percent of returning males after a season with 30 breeding pairs.
- (b) What percent of the year-to-year variation in percent of returning males is accounted for by the straightline relationship with number of breeding pairs the previous year?

Other types of software displays, including Minitab, Fathom, and applet screen captures, appear throughout the book to help you learn to read and interpret many different kinds of output.

Overview

What Is Statistics?

Does listening to music while studying help or hinder learning? If an athlete fails a drug test, how sure can we be that she took a banned substance? Does having a pet help people live longer? How well do SAT scores predict college success? Do most people recycle? Which of two diets will help obese children lose more weight and keep it off? Should a poker player go "all in" with pocket aces? Can a new drug help people quit smoking? How strong is the evidence for global warming?

These are just a few of the questions that statistics can help answer. But what is statistics? And why should you study it?

Statistics Is the Science of Learning from Data

Data are usually numbers, but they are not "just numbers." Data are numbers with a context. The number 10.5, for example, carries no information by itself. But if we hear that a family friend's new baby weighed 10.5 pounds at birth, we congratulate her on the healthy size of the child. The context engages our knowledge



about the world and allows us to make judgments. We know that a baby weighing 10.5 pounds is quite large, and that a human baby is unlikely to weigh 10.5 ounces or 10.5 kilograms. The context makes the number meaningful.

In your lifetime, you will be bombarded with data and statistical information. Poll results, television ratings, music sales, gas prices, unemployment rates, medical study outcomes, and standardized test scores are discussed daily in the media. Using data effectively is a large and growing part of most professions. A solid understanding of statistics will enable you to make sound, data-based decisions in your career and everyday life.

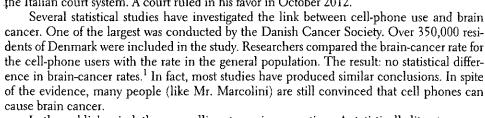
Data Beat Personal Experiences

It is tempting to base conclusions on your own experiences or the experiences of those you know. But our experiences may not be typical. In fact, the incidents that stick in our memory are often the unusual ones.

Do cell phones cause brain cancer?

Italian businessman Innocente Marcolini developed a brain tumor at age 60. He also talked on a cellular phone up to 6 hours per day for 12 years as part of his job. Mr. Marcolini's physician suggested that the

brain tumor may have been caused by cell-phone use. So Mr. Marcolini decided to file suit in the Italian court system. A court ruled in his favor in October 2012.



In the public's mind, the compelling story wins every time. A statistically literate person knows better. Data are more reliable than personal experiences because they systematically describe an overall picture rather than focus on a few incidents.



Where the Data Come from Matters



Are you kidding me?

The famous advice columnist Ann Landers once asked her readers, "If you had it to do over again, would you have children?" A few weeks later, her column was headlined "70% OF PARENTS SAY KIDS NOT WORTH IT." Indeed, 70% of the nearly 10,000 parents who wrote in said they would not have children if they could make the choice again. Do you believe that 70% of all parents regret having children?

You shouldn't. The people who took the trouble to write Ann Landers are not representative of all parents. Their letters showed that many of them were angry with their children. All we know from these data is that there are some unhappy parents out there. A statistically designed poll, unlike Ann Landers's appeal, targets specific people chosen in a way that gives all parents the same chance to be asked. Such a poll showed that 91% of parents would have children again.

Where data come from matters a lot. If you are careless about how you get your data, you may announce 70% "No" when the truth is close to 90% "Yes."

Who talks more—women or men?

According to Louann Brizendine, author of *The Female Brain*, women say nearly three times as many words per day as men. Skeptical researchers devised a study to test this claim. They used electronic devices to record the talking patterns of 396 university students from Texas, Arizona, and Mexico. The device was programmed to record 30 seconds of sound every 12.5 minutes without the carrier's knowledge. What were the results?

According to a published report of the study in Scientific American, "Men showed a slightly wider variability in words uttered. . . . But in the end, the sexes came out just about even in the daily averages: women at 16,215 words and men at 15,669." When asked where she got her figures, Brizendine admitted that she used unreliable sources.

The most important information about any statistical study is how the data were produced. Only carefully designed studies produce results that can be trusted.

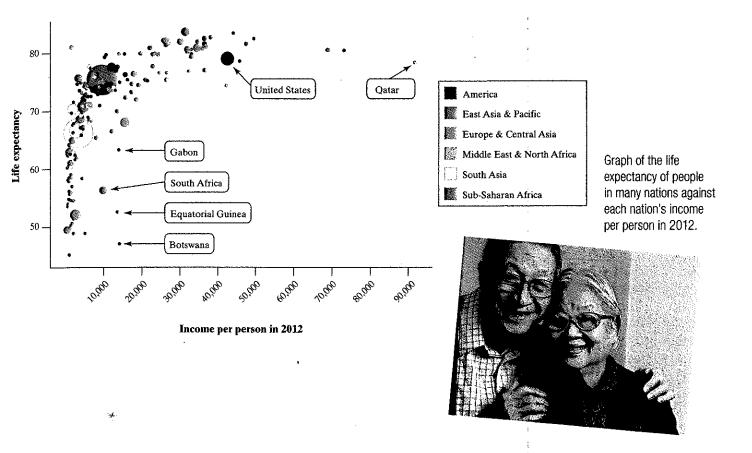
Always Plot Your Data

Yogi Berra, a famous New York Yankees baseball player known for his unusual quotes, had this to say: "You can observe a lot just by watching." That's a motto for learning from data. A carefully chosen graph is often more instructive than a bunch of numbers.

Do people live longer in wealthier countries?

The Gapminder Web site, www.gapminder.org, provides loads of data on the health and well-being of the world's inhabitants. The graph on the next pages displays some data from Gapminder.⁴ The individual points represent all the world's nations for which data are available. Each point shows the income per person and life expectancy in years for one country.

We expect people in richer countries to live longer. The overall pattern of the graph does show this, but the relationship has an interesting shape. Life expectancy rises very quickly as personal income increases and then levels off. People in very rich countries like the United States live no longer than people in poorer but not extremely poor nations. In some less wealthy countries, people live longer than in the United States. Several other nations stand out in the graph. What's special about each of these countries?



Variation Is Everywhere

Individuals vary. Repeated measurements on the same individual vary. Chance outcomes—like spins of a roulette wheel or tosses of a coin—vary. Almost everything varies over time. Statistics provides tools for understanding variation.

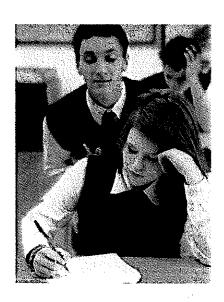
Have most students cheated on a test?

Researchers from the Josephson Institute were determined to find out. So they surveyed about 23,000 students from 100 randomly selected schools (both public and private) nationwide. The question they asked was "How many times have you cheated during a test at school in the past year?" Fifty-one percent said they had cheated at least once.⁵

If the researchers had asked the same question of *all* high school students, would exactly 51% have answered "Yes"? Probably not. If the Josephson Institute had selected a different sample of about 23,000 students to respond to the survey, they would probably have gotten a different estimate. *Variation is everywhere!*

Fortunately, statistics provides a description of how the sample results will vary in relation to the actual population percent. Based on the sampling method that this study used, we can say that the estimate of 51% is very likely to be within 1% of the true population value. That is, we can be quite confident that between 50% and 52% of all high school students would say that they have cheated on a test.

Because variation is everywhere, conclusions are uncertain. Statistics gives us a language for talking about uncertainty that is understood by statistically literate people everywhere.



Designing Studies

case study

Can Magnets Help Reduce Pain?

Early research showed that magnetic fields affected living tissue in humans. Some doctors have begun to use magnets to treat patients with chronic pain. Scientists wondered whether this type of therapy really worked. They designed a study to find out.

Fifty patients with chronic pain were recruited for the study. A doctor identified a painful site on each patient and asked him or her to rate the pain on a scale from 0 (mild pain) to 10 (severe pain). Then, the doctor selected a sealed envelope containing a magnet at random from a box with a mixture of active and inactive magnets. That way, neither the doctor nor the patient knew which type of magnet was being used. The chosen magnet was applied to the site of the pain for 45 minutes. After "treatment," each patient was again asked to rate the level of pain from 0 to 10.

In all, 29 patients were given active magnets and 21 patients received inactive magnets. Scientists decided to focus on the improvement in patients' pain ratings. Here they are, grouped by the type of magnet used:¹

Active: 10 6 1 10 6 8 5 5 6 8 7 8 7 6 4 4 7 10 6 10 6 5 5 1 0 0 0 0 1

Inactive: 4 3 5 2 1 4 1 0 0 1 0 0 0 0 0 0 0 1

What do the data tell us about whether the active magnets helped reduce pain? By the end of the chapter, you'll be ready to interpret the results of this study.



Introduction

You can hardly go a day without hearing the results of a statistical study. Here are some examples:

 The National Highway Traffic Safety Administration (NHTSA) reports that seat belt use in passenger vehicles increased from 84% in 2011 to 86% in 2012.²



- According to a recent survey, U.S. teens aged 13 to 18 spend an average of 26.8 hours per week online. Although 59% of the teens said that posting personal information or photos online is unsafe, 62% said they had posted photos of themselves.³
- A recent study suggests that lack of sleep increases the risk of catching a cold.⁴
- For their final project, two AP® Statistics students showed that listening to music while studying decreased subjects' performance on a memory task.⁵

Can we trust these results? As you'll learn in this chapter, that depends on how the data were produced. Let's take a closer look at where the data came from in each of these studies.

Each year, the NHTSA conducts an observational study of seat belt use in vehicles. The NHTSA sends trained observers to record the actual behavior of people in vehicles at randomly selected locations across the country. The idea of an observational study is simple: you can learn a lot just by watching. Or by asking a few questions, as in the survey of teens' online habits. Harris Interactive conducted this survey using a "representative sample" of 655 U.S. 13- to 18-year-olds. Both of these studies use information from a sample to draw conclusions about some larger population. Section 4.1 examines the issues involved in sampling and surveys.

In the sleep and catching a cold study, 153 volunteers took part. They answered questions about their sleep habits over a two-week period. Then, researchers gave them a virus and waited to see who developed a cold. This was a complicated observational study. Compare this with the *experiment* performed by the AP® Statistics students. They recruited 30 students and divided them into two groups of 15 by drawing names from a hat. Students in one group tried to memorize a list of words while listening to music. Students in the other group tried to memorize the same list of words while sitting in silence. Section 4.2 focuses on designing experiments.

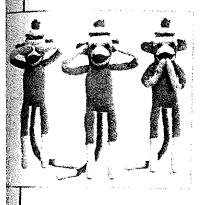
The goal of many statistical studies is to show that changes in one variable cause changes in another variable. In Section 4.3, we'll look at why establishing causation is so difficult, especially in observational studies. We'll also consider some of the ethical issues involved in planning and conducting a study.

Here's an Activity that gives you a preview of what lies ahead.

ACTIVITY | See no evil, hear no evil?

MATERIALS:

Two index cards, each with 10 distinct numbers from 00 to 99 written on it (prepared by your teacher); clock, watch, or stopwatch to measure 30 seconds; and a coin for each pair of students



Confucius said, "I hear and I forget. I see and I remember. I do and I understand." Do people really remember what they see better than what they hear? In this Activity, you will perform an experiment to try to find out.

- 1. Divide the class into pairs of students by drawing names from a hat.
- 2. Your teacher will give each pair two index cards with 10 distinct numbers from 00 to 99 on them. Do not look at the numbers until it is time for you to do the experiment.
- 3. Flip a coin to decide which of you is Student 1 and which is Student 2. Shuffle the index cards and deal one face down to each partner.
- 4. Student 1 will be the first to attempt a memory task while Student 2 keeps time.

Directions: Study the numbers on the index card for 30 seconds. Then turn the card over. Recite the alphabet aloud (A, B, C, and so on). Then tell your partner what you think the numbers on the card are. You may not say more than 10 numbers! Student 2 will record how many numbers you recalled correctly.

- 5. Now it's Student 2's turn to do a memory task while Student 1 records the data. Directions: Your partner will read the numbers on your index card aloud three times slowly. Next, you will recite the alphabet aloud (A, B, C, and so on) and then tell your partner what you think the numbers on the card are. You may not say more than 10 numbers! Student 1 will record how many numbers you recalled correctly.
- 6. Your teacher will scale and label axes on the board for parallel dotplots of the results. Plot how many numbers you remembered correctly on the appropriate graph.
- 7. Did students in your class remember numbers better when they saw them or when they heard them? Give appropriate evidence to support your answer.
- 8. Based on the results of this experiment, can we conclude that people in general remember better when they see than when they hear? Why or why not?

4.1

Sampling and Surveys

WHAT YOU WILL LEARN

By the end of the section, you should be able to:

- Identify the population and sample in a statistical study.
- Identify voluntary response samples and convenience samples. Explain how these sampling methods can lead to bias.
- Describe how to obtain a random sample using slips of paper, technology, or a table of random digits.
- Distinguish a simple random sample from a stratified random sample or cluster sample. Give the advantages and disadvantages of each sampling method.
- Explain how undercoverage, nonresponse, question wording, and other aspects of a sample survey can lead to bias.

Suppose we want to find out what percent of young drivers in the United States text while driving. To answer the question, we will survey 16- to 20-year-olds who live in the United States and drive. Ideally, we would ask them all (take a census). But contacting every driver in this age group wouldn't be practical: it would take too much time and cost too much money. Instead, we put the question to a sample chosen to represent the entire population of young drivers.

DEFINITION: Population, census, and sample

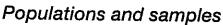
The population in a statistical study is the entire group of individuals we want information about. A census collects data from every individual in the population.

A sample is a subset of individuals in the population from which we actually collect data.

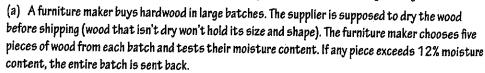
The distinction between population and sample is basic to statistics. To make sense of any sample result, you must know what population the sample represents. Here's an example that illustrates this distinction and also introduces some major uses of sampling.







PROBLEM: Identify the population and the sample in each of the following settings.

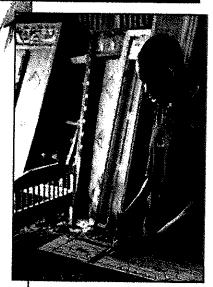


(b) Each week, the Gallup Poll questions a sample of about 1500 adult U.S. residents to determine national opinion on a wide variety of issues.

SOLUTION:

- (a) The population is all the pieces of hardwood in a batch. The sample is the five pieces of wood that are selected from that batch and tested for moisture content.
- (b) Gallup's population is all adult U.S. residents. Their sample is the 1500 adults who actually respond to the survey questions.

For Practice Try Exercise |



The Idea of a Sample Survey

We often draw conclusions about a whole population on the basis of a sample. Have you ever tasted a sample of ice cream and ordered a cone if the sample tastes good? Because ice cream is fairly uniform, the single taste represents the whole. Choosing a representative sample from a large and varied population (like all young U.S. drivers) is not so easy. The first step in planning a sample survey is to say exactly what population we want to describe. The second step is to say exactly what we want to measure, that is, to give exact definitions of our variables.

We reserve the term "sample survey" for studies that use an organized plan to choose a sample that represents some specific population, like the pieces of hardwood and the U.S. adults in the previous example. By our definition, the population in a sample survey can consist of people, animals, or things. Some people use the terms "survey" or "sample survey" to refer only to studies in which people are asked one or more questions, like the Gallup Poll of the last example. We'll avoid this more restrictive terminology.



How Does the Current Population Survey Work?

A sample survey

One of the most important government sample surveys in the United States is the monthly Current Population Survey (CPS). The CPS contacts about 60,000 households each month. It produces the monthly unemployment rate and lots of other economic and social information. To measure unemployment, we must first specify the population we want to describe. The CPS defines its population as all U.S. residents (legal or not) 16 years of age and over who are civilians and are not in an institution such as a prison. The unemployment rate announced in the news refers to this specific population.

What does it mean to be "unemployed"? Someone who is not looking for work—for example, a full-time student—should not be called unemployed just because she is not working for pay. If you are chosen for the CPS sample, the interviewer first asks whether you are available to work and whether you actually looked for work in the past four weeks. If not, you are neither employed nor unemployed—you are not in the labor force.

If you are in the labor force, the interviewer goes on to ask about employment. If you did any work for pay or in your own business during the week of the survey, you are employed. If you worked at least 15 hours in a family business without pay, you are employed. You are also employed if you have a job but didn't work because of vacation, being on strike, or other good reason. An unemployment rate of 9.7% means that 9.7% of the sample was unemployed, using the exact CPS definitions of both "labor force" and "unemployed."

The final step in planning a sample survey is to decide how to choose a sample from the population. Let's take a closer look at some good and not-so-good sampling methods.

How to Sample Badly

Suppose we want to know how long students at a large high school spent doing homework last week. We might go to the school library and ask the first 30 students we see about their homework time. The sample we get is known as a convenience sample.

The sampling method that yields a convenience sample is called convenience sampling. Other sampling methods are named in similarly obvious ways!

DEFINITION: Convenience sample

Choosing individuals from the population who are easy to reach results in a convenience sample.

Convenience sampling often produces unrepresentative data. Consider our sample of 30 students from the school library. It's unlikely that this convenience sample accurately represents the homework habits of all students at the high school. In fact, if we were to repeat this sampling process again and again, we would almost always overestimate the average homework time in the population. Why? Because students who hang out in the library tend to be more studious. This is bias: using a method that favors some outcomes over others.

DEFINITION: Bias

The design of a statistical study shows **bias** if it would consistently underestimate or consistently overestimate the value you want to know.

AP® EXAM TIP If you're asked to describe how the design of a study leads to bias, you're expected to do two things: (1) identify a problem with the design, and (2) explain how this problem would lead to an underestimate or overestimate. Suppose you were asked, "Explain how using your statistics class as a sample to estimate the proportion of all high school students who own a graphing calculator could result in bias." You might respond, "This is a convenience sample. It would probably include a much higher proportion of students with a graphing calculator than in the population at large because a graphing calculator is required for the statistics class. So this method would probably lead to an overestimate of the actual population proportion."

Bias is not just bad luck in one sample. It's the result of a bad study design that will consistently miss the truth about the population in the same way. Convenience samples are almost guaranteed to show bias. So are voluntary response samples.



Voluntary response samples are also known as self-selected samples.

DEFINITION: Voluntary response sample

A **voluntary response sample** consists of people who choose themselves by responding to a general invitation.

Call-in, text-in, write-in, and many Internet polls rely on voluntary response samples. People who choose to participate in such surveys are usually not representative of some larger population of interest. Voluntary response samples attract people who feel strongly about an issue, and who often share the same opinion. That leads to bias.



The Internet brings voluntary response samples to the computer nearest you. Visit www.misterpoll.com to become part of the sample in any of dozens of online polls. As the site says, "None of these polls are 'scientific,' but do represent the collective opinion of everyone who participates." Unfortunately, such polls don't tell you anything about the views of the population.





Illegal Immigration

Online polls

Former CNN commentator Lou Dobbs doesn't like illegal immigration. One of his shows was largely devoted to attacking a proposal to offer driver's licenses to illegal immigrants. During the show, Mr. Dobbs invited his viewers to go to loudobbs.com to vote on the question "Would you be more or less likely to vote for a presidential candidate who supports giving driver's licenses to illegal aliens? The result: 97% of the 7350 people who voted by the end of the show said, "Less likely."

PROBLEM: What type of sample did Mr. Dobbs use in his poll? Explain how this sampling method could lead to bias in the poll results.

SOLUTION: Mr. Dobbs used a voluntary response sample: people chose to go online and respond. Those who voted were viewers of Mr. Dobbs's program, which means that they are likely to support his views. The 97% poli result is probably an extreme overestimate of the percent of people in the population who would be less likely to support a presidential candidate with this position.

For Practice Try Exercise





CHECK YOUR UNDERSTANDING

For each of the following situations, identify the sampling method used. Then explain how the sampling method could lead to bias.

- 1. A farmer brings a juice company several crates of oranges each week. A company inspector looks at 10 oranges from the top of each crate before deciding whether to buy all the oranges.
- The ABC program Nightline once asked whether the United Nations should continue to have its headquarters in the United States. Viewers were invited to call one telephone number to respond "Yes" and another for "No." There was a charge for calling either number. More than 186,000 callers responded, and 67% said "No."

How to Sample Well: Simple Random Sampling

In convenience sampling, the researcher chooses easy-to-reach members of the population. In voluntary response sampling, people decide whether to join the sample. Both sampling methods suffer from bias due to personal choice. The best way to avoid this problem is to let chance choose the sample. That's the idea of random sampling.

DEFINITION: Random Sampling

Random sampling involves using a chance process to determine which members of a population are included in the sample.

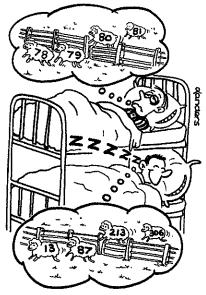
In everyday life, some people use the word "random" to mean haphazard, as in "that's so random." In statistics, random means "due to chance."

Don't say that a sample was chosen at random if a chance process wasn't used to select the individuals.

The easiest way to choose a random sample of n people is to write their names on identical slips of paper, put the slips in a hat, mix them well, and pull out slips one at a time until you have n of them. An alternative would be to give each member of the population a distinct number and to use the "hat method" with these numbers instead of people's names. Note that this version would work just as well if the population consisted of animals or things. The resulting sample is called a simple random sample, or SRS for short.

DEFINITION: Simple Random Sample (SRS)

A **simple random sample** (SRS) of size n is chosen in such a way that every group of n individuals in the population has an equal chance to be selected as the sample.



Statisticians Fall asleep Faster by taking a random sample of sheep.

An SRS gives every possible sample of the desired size an equal chance to be chosen. It also gives each member of the population an equal chance to be included in the sample. Picture drawing 20 slips (the sample) from a hat containing 200 identical slips (the population). Any 20 slips have the same chance as any other 20 to be chosen. Also, each slip has a 1-in-10 chance (20/200) of being selected.

Some other random sampling methods give each member of the population, but not each sample, an equal chance. We'll look at some of these later.

How to Choose a Simple Random Sample The hat method won't work well if the population is large. Imagine trying to take a simple random sample of 1000 U.S. adults! In practice, most people use random numbers generated by technology to choose samples.

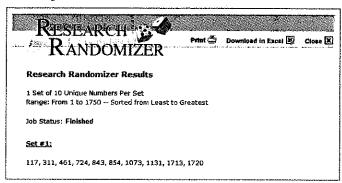


Teens on the Internet

Choosing an SRS with technology

The principal at Canyon del Oro High School in Arizona wants student input about limiting access to certain Internet sites on the school's computers. He asks the AP® Statistics teacher, Mr. Tabor, to select a "representative sample" of 10 students. Mr. Tabor decides to take an SRS from the 1750 students enrolled this year.

He gets an alphabetical roster from the registrar's office, and numbers the students from 1 to 1750. Then Mr. Tabor uses the random number generator at www.randomizer.org to choose 10 distinct numbers between 1 and 1750:



The 10 students on the roster that correspond to the chosen numbers will be on the principal's committee.

This example highlights the steps in choosing a simple random sample with technology.

CHOOSING AN SRS WITH TECHNOLOGY

It is standard practice to use *n* for the sample size and *N* for the population size

Step 1: Label. Give each individual in the population a distinct numerical label from 1 to *N*.

Step 2: Randomize. Use a random number generator to obtain n different integers from 1 to N.

You can also use a graphing calculator to choose an SRS.



10. TECHNOLOGY CORNER

CHOOSING AN SRS

TI-Nspire instructions in Appendix B; HP Prime instructions on the book's Web site.

Let's use a graphing calculator to select an SRS of 10 students from the Canyon del Oro High School roster.

1. Check that your calculator's random number generator is working properly.

TI-83/84

Press MATH, then select PRB and randInt (. Complete the command randInt (1,1750) and press ENTER.

TI-89

• Press CATALOG, then F3 (Flash Apps) and choose randInt (. Complete the command TIStat. randInt (1, 1750) and press ENTER.

Compare your results with those of your classmates. If several students got the same number, you'll need to seed your calculator's random integer generator with different numbers before you proceed. Directions for doing this are given in the Annotated Teacher's Edition.

Randomly generate 10 distinct numbers from 1 to 1750.
 Do randInt (1, 1750) again. Keep pressing ENTER until you have chosen 10 different labels.

| NORMAL FLOAT AUTO REAL RADIA! | 1 (1) |
|-------------------------------|--|
| randInt(1,1750) | 50 X 4 - 10 M 50 |
| randInt(1,1750) | 139 |
| | 1126 |
| randInt(1,1750) | 920 |
| randInt(1.1750) | |
| | 1089 |
| <u></u> | |

| F1- F2- F3- F4- F5 Too is A19ebra Calc Other Prymin | F6- |
|--|----------------|
| ■tistat.randint(1, | 1750) |
| <pre>* tistat.randint(1,</pre> | 1125. 1750) |
| *tistat.randint(1, | 1265. |
| | 240. |
| TIStat.randInt(1,17 | (C 3/30 |

Note: If you have a TI-83/84 with OS 2.55 or later, you can use the command RandIntNoRep (1, 1750) to sort the numbers from 1 to 1750 in random order. The first 10 numbers listed give the labels of the chosen students.

If you don't have technology handy, you can use a table of random digits to choose an SRS. We have provided a table of random digits at the back of the book (Table D). Here is an excerpt.

| | | | Table | D Rando | m digits | | ⊈ € | |
|------|-------|-------|-------|---------|----------|-----------------------------|------------|-------|
| LINE | | | | | | eres represent the particle | | |
| 101 | 19223 | 95034 | 05756 | 28713 | 96409 | 12531 | 42544 | 82853 |
| 102 | 73676 | 47150 | 99400 | 01927 | 27754 | 42648 | 82425 | 36290 |
| 103 | 45467 | 71709 | 77558 | 00095 | 32863 | 29485 | 82226 | 90056 |

You can think of this table as the result of someone putting the digits 0 to 9 in a hat, mixing, drawing one, replacing it, mixing again, drawing another, and so on. The digits have been arranged in groups of five within numbered rows to make the table easier to read. The groups and rows have no special meaning—Table D is just a long list of randomly chosen digits. As with technology, there are two steps in using Table D to choose a random sample.

HOW TO CHOOSE AN SRS USING TABLE D

Step 1: Label. Give each member of the population a numerical label with the same number of digits. Use as few digits as possible.

Step 2: Randomize. Read consecutive groups of digits of the appropriate length from left to right across a line in Table D. Ignore any group of digits that wasn't used as a label or that duplicates a label already in the sample. Stop when you have chosen *n* different labels.

Your sample contains the individuals whose labels you find.

Always use the shortest labels that will cover your population. For instance, you can label up to 100 individuals with two digits: 01, 02, . . . , 99, 00. As standard practice, we recommend that you begin with label 1 (or 01 or 001 or 0001, as needed). Reading groups of digits from the table gives all individuals the same chance to be chosen because all labels of the same length have the same chance

to be found in the table. For example, any pair of digits in the table is equally likely to be any of the 100 possible labels 01, 02, ..., 99, 00. Here's an example that shows how this process works.



Spring Break!

Choosing an SRS with Table D



The school newspaper is planning an article on family-friendly places to stay over spring break at a nearby beach town. The editors intend to call 4 randomly chosen hotels to ask about their amenities for families with children. They have an alphabetized list of all 28 hotels in the town.

PROBLEM: Use Table D at line 130 to choose an SRS of 4 hotels for the editors to call.

SOLUTION: We'll use the two-step process for selecting an SRS using Table D.

Step 1: Label. Two digits are needed to label the 28 hotels. We have added labels 01 to 28 to the alphabetized list of hotels below.

| 01 Aloha Kai | 08 Captiva | 15 Palm Tree | 22 Sea Shell |
|-----------------|-----------------|---------------|--------------------|
| 02 Anchor Down | 09 Casa del Mar | 16 Radisson | 23 Silver Beach |
| 03 Banana Bay | 10 Coconuts | 17 Ramada | 24 Sunset Beach |
| 04 Banyan Tree | 11 Diplomat | 18 Sandpiper | 25 Tradewinds |
| 05 Beach Castle | 12 Holiday Inn | 19 Sea Castle | 26 Tropical Breeze |
| 06 Best Western | 13 Lime Tree | 20 Sea Club | 27 Tropical Shores |
| 07 Cabana | 14 Outrigger | 21 Sea Grape | 28 Veranda |

Step 2: Randomize. To use Table D, start at the left-hand side of line 130 and read two-digit groups. Skip any groups that aren't between 01 and 28, as well as any repeated groups. Continue until you have chosen four hotels. Here is the beginning of line 130:

| 69051 | 6481 | 17 <i>8</i> | 7174 | 09517 | 84534 | 064 | 89 8 | 7201 | 97245 | |
|-----------|------------|-------------|---------|-------|---------|--------|---------|---------|--------|--|
| The first | 10 two-dig | jit group | s are | | | | | | | |
| 69 | 05 | 16 | 48 | 17 | 87 | 17 | 40 | 95 | 17 | |
| Skip | 1 | 1 | Skip | 1 | Skip | Skip | 5kip | Skip | Skip | |
| Too big | | | Too big | | Too big | Repeat | Too big | Too bia | Repeat | |

We skip 5 of these 10 groups because they are too high (over 28) and 2 because they are repeats (both 17s). The hotels labeled 05, 16, and 17 go into the sample. We need one more hotel to complete the sample. Continuing along line 130:

| 84 | 53 | 40 | 64 | 89 | 87 | 20 |
|---------|---------|---------|---------|---------|---------|----------|
| Skip | Skip | Skip | Skip | Skip | Skip | <i>∠</i> |
| Too big | Too big | Too big | Too big | Too biq | Too bia | |

Our SRS of 4 hotels for the editors to contact is 05 Beach Castle, 16 Radisson, 17 Ramada, and 20 Sea Club.

We can trust results from an SRS, as well as from other types of random samples that we will meet later, because the use of impersonal chance avoids bias. The following activity shows why random sampling is so important.

ACTIVITY Who Wrote the Federalist Papers?

The Federalist Papers are a series of 85 essays supporting the ratification of the U.S. Constitution. At the time they were published, the identity of the authors was a secret known to just a few people. Over time, however, the authors were identified as Alexander Hamilton, James Madison, and John Jay. The authorship of 73 of the essays is fairly certain, leaving 12 in dispute. However, thanks in some part to statistical analysis, most scholars now believe that the 12 disputed essays were written by Madison alone or in collaboration with Hamilton.

There are several ways to use statistics to help determine the authorship of a disputed text. One method is to estimate the average word length in a disputed text and compare it to the average word lengths of works where the authorship is not in dispute.

The following passage is the opening paragraph of Federalist Paper #51,9 one of the disputed essays. The theme of this essay is the separation of powers between the three branches of government.

To what expedient, then, shall we finally resort, for maintaining in practice the necessary partition of power among the several departments, as laid down in the Constitution? The only answer that can be given is, that as all these exterior provisions are found to be inadequate, the defect must be supplied, by so contriving the interior structure of the government as that its several constituent parts may, by their mutual relations, be the means of keeping each other in their proper places. Without presuming to undertake a full development of this important idea, I will hazard a few general observations, which may perhaps place it in a clearer light, and enable us to form a more correct judgment of the principles and structure of the government planned by the convention.

- 1. Choose 5 words from this passage. Count the number of letters in each of the words you selected, and find the average word length.
- 2. Your teacher will draw and label a horizontal axis for a class dotplot. Plot the average word length you obtained in Step 1 on the graph.
- 3. Use a table of random digits or a random number generator to select a simple random sample of 5 words from the 130 words in the opening passage. Count the number of letters in each of the words you selected, and find the average word length.
- 4. Your teacher will draw and label another horizontal axis with the same scale for a comparative class dotplot. Plot the average word length you obtained in Step 3 on the graph.
- 5. How do the dotplots compare? Can you think of any reasons why they might be different? Discuss with your classmates.

Other Random Sampling Methods

The basic idea of sampling is straightforward: take an SRS from the population and use your sample results to gain information about the population. Unfortunately, it's usually difficult to get an SRS from the population of interest. Imagine trying to get a simple random sample of all the batteries produced in one day at a factory. Or an SRS of all U.S. high school students. In either case, it would be difficult to obtain an accurate list of the population from which to draw the sample. It would also be very time-consuming to collect data from each individual that's randomly selected. Sometimes, there are also statistical advantages to using more complex sampling methods.

One of the most common alternatives to an SRS involves sampling groups (strata) of similar individuals within the population separately. Then these separate "subsamples" are combined to form one stratified random sample.

Stratum is singular. Strata are plural.

DEFINITION: Stratified random sample and strata

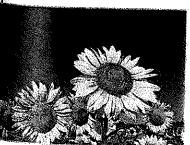
To get a **stratified random sample**, start by classifying the population into groups of similar individuals, called **strata**. Then choose a separate SRS in each stratum and combine these SRSs to form the sample.

Choose the strata based on facts known before the sample is taken. For example, in a study of sleep habits on school nights, the population of students in a large high school might be divided into freshman, sophomore, junior, and senior strata. In a preelection poll, a population of election districts might be divided into urban, suburban, and rural strata. Stratified random sampling works best when the individuals within each stratum are similar with respect to what is being measured and when there are large differences between strata. The following Activity makes this point clear.

ACTIVITY Sampling sunflowers

MATERIALS:

Calculator for each student



A British farmer grows sunflowers for making sunflower oil. Her field is arranged in a grid pattern, with 10 rows and 10 columns as shown in the figure on the next page. Irrigation ditches run along the top and bottom of the field. The farmer would like to estimate the number of healthy plants in the field so she can project

how much money she'll make from selling them. It would take too much time to count the plants in all 100 squares, so she'll accept an estimate based on a sample of 10 squares.

- I. Use Table D or technology to take a simple random sample of 10 grid squares. Record the location (for example, B6) of each square you select.
- 2. This time, you'll take a stratified random sample using the *rows* as strata. Use Table D or technology to randomly select one square from each (horizontal) row. Record the location of each square—for example, Row 1: G, Row 2: B, and so on.

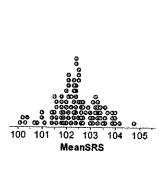
| | Α | В | С | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|--------|---|---|---|
| 1 | | | | | | | | | | |
| 2 | | | | | | | | | | |
| 3 | | | | | | | | | | |
| 4 | | | | | | | | | | |
| 5 | | | | | | | | | | |
| 6 | | | | | | | | | | - |
| 7 | | | | | | | | | | |
| 8 | | | | | | | ****** | | | |
| 9 | | | | | | | | | | |

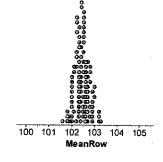
- 3. Now, take a stratified random sample using the *columns* as strata. Use Table D or technology to randomly select one square from each (vertical) column. Record the location of each square—for example, Column A: 4, Column B: 1, and so on.
- 4. The table on page N/DS-5 in the back of the book gives the actual number of sunflowers in each grid square. Use the information provided to calculate your estimate of the mean number of sunflowers per square for each of your samples in Steps 1, 2, and 3.
- 5. Make comparative dotplots showing the mean number of sunflowers obtained using the three different sampling methods for all members of the class. Describe any similarities and differences you see.
- 6. Your teacher will provide you with the mean number of sunflowers in the population of all 100 grid squares in the field. How did the three sampling methods do?

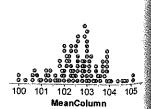
The dotplots below show the mean number of healthy plants in 100 samples using each of the three sampling methods in the Activity: simple random sampling, stratified random sampling with rows of the field as strata, and stratified random sampling with columns of the field as strata. Notice that all three distributions are centered at about 102.5, the true mean number of healthy plants in all squares of the field. That makes sense because random sampling yields accurate estimates of unknown population values.

One other detail stands out in the graphs. There is much less variability in the estimates using stratified random sampling with the rows as strata. The table on page N/DS-5 shows the actual number of healthy sunflowers in each grid square. Notice that the squares within each row contain a similar number of healthy plants but there are big differences between rows. When we can choose strata that are "similar within but different between," stratified random samples give more precise estimates than simple random samples of the same size.

Why didn't using the columns as strata reduce the variability of the estimates in a similar way? Because the squares within each column have very different numbers of healthy plants.







Both simple random sampling and stratified random sampling are hard to use when populations are large and spread out over a wide area. In that situation, we'd

prefer a method that selects groups (clusters) of individuals that are "near" one another. That's the idea of a cluster sample.

DEFINITION: Cluster sample and clusters

To get a **cluster sample**, start by classifying the population into groups of individuals that are located near each other, called **clusters**. Then choose an SRS of the clusters. All individuals in the chosen clusters are included in the sample.

In a cluster sample, some people take an SRS from each cluster rather than including all members of the cluster.

Cluster samples are often used for practical reasons, like saving time and money. Cluster sampling works best when the clusters look just like the population but on a smaller scale. Imagine a large high school that assigns its students to homerooms alphabetically by last name. The school administration is considering a new schedule and would like student input. Administrators decide to survey 200 randomly selected students. It would be difficult to track down an SRS of 200 students, so the administration opts for a cluster sample of homerooms. The principal (who knows some statistics) takes a simple random sample of 8 homerooms and gives the survey to all 25 students in each homeroom.

Cluster samples don't offer the statistical advantage of better information about the population that stratified random samples do. That's because clusters are often chosen for ease so they may have as much variability as the population itself

itself.

Be sure you understand the difference between strata and clusters. We want each stratum to contain similar individuals and for there to be large differences between strata. For a cluster sample, we'd like each cluster to look just like the population, but on a smaller scale. Here's an example that compares the random sampling methods we have discussed so far.

Remember: strata are ideally "similar within, but different between," while dusters are ideally "different within, but similar between."



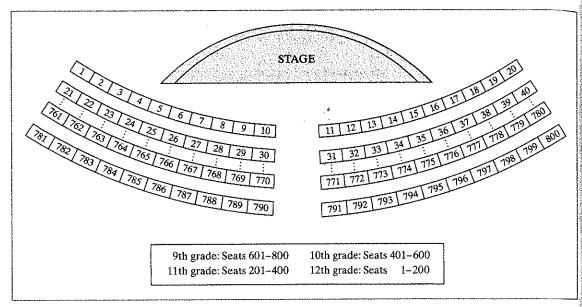
Sampling at a School Assembly

Strata or clusters?



The student council wants to conduct a survey during the first five minutes of an all-school assembly in the auditorium about use of the school library. They would like to announce the results of the survey at the end of the assembly. The student council president asks your statistics class to help carry out the survey.

PROBLEM: There are 800 students present at the assembly. A map of the auditorium is shown on the next page. Note that students are seated by grade level and that the seats are numbered from 1 to 800.



Describe how you would use your calculator to select 80 students to complete the survey with each of the following:

- (a) Simple random sample
- (b) Stratified random sample
- (c) Cluster sample

SOLUTION:

- (a) To take an SRS, we need to choose 80 of the seat numbers at random. Use randint (1,800) on your calculator until 80 different seats are selected. Then give the survey to the students in those seats.
- (b) The students in the assembly are seated by grade level. Because students' library use might be similar within grade levels but different across grade levels, we'll use the grade level seating areas as our strata. Within each grade's seating area, we'll select 20 seats at random. For the 9th grade, use randint(601,800) to select 20 different seats. Use randint(401,600) to pick 20 different sophomore seats, randint(201,400) to get 20 different junior seats, and randint(1,200) to choose 20 different senior seats. Give the survey to the students in the selected seats.
- (c) With the way students are seated, each column of seats from the stage to the back of the auditorium could be used as a cluster. Note that each cluster contains students from all four grade levels, so each should represent the population well. Because there are 20 clusters, each with 40 seats, we need to choose 2 clusters at random to get 80 students for the survey. Use randint(1,20) to select two clusters, and then give the survey to all 40 students in each column of seats.

Note that cluster sampling is much more efficient than finding 80 seats scattered about the auditorium, as required by both of the other sampling methods.

For Practice Try Exercise

Most large-scale sample surveys use *multistage samples* that combine two or more sampling methods. For example, the U.S. Census Bureau carries out a monthly Current Population Survey (CPS) of about 60,000 households. Researchers start by choosing a stratified random sample of neighborhoods in 756 of the 2007 geographical areas in the United States. Then they divide each neighborhood into clusters of four nearby households and select a cluster sample to interview.

Analyzing data from sampling methods more complex than an SRS takes us beyond basic statistics. But the SRS is the building block of more elaborate methods, and the principles of analysis remain much the same for these other methods.



The manager of a sports arena wants to learn more about the financial status of the people who are attending an NBA basketball game. He would like to give a survey to a representa-

tive sample of the more than 20,000 fans in attendance. Ticket prices for the game vary a great deal: seats near the court cost over \$100 each, while seats in the top rows of the arena cost \$25 each. The arena is divided into 30 numbered sections, from 101 to 130. Each section has rows of seats labeled with letters from A (nearest the court) to ZZ (top row of the arena).

- Explain why it might be difficult to give the survey to an SRS of 200 fans.
- Which would be a better way to take a stratified random sample of fans: using the lettered rows or the numbered sections as strata? Explain.
- Which would be a better way to take a cluster sample of fans: using the lettered rows or the numbered sections as clusters? Explain.

Inference for Sampling

The purpose of a sample is to give us information about a larger population. The process of drawing conclusions about a population on the basis of sample data is called inference because we infer information about the population from what we know about the sample.

Inference from convenience samples or voluntary response samples would be misleading because these methods of choosing a sample are biased. We are almost certain that the sample does not fairly represent the population. The first reason to rely on random sampling is to avoid bias in choosing a sample.

Still, it is, unlikely that results from a random sample are exactly the same as for the entire population. Sample results, like the unemployment rate obtained from the monthly Current Population Survey, are only estimates of the truth about the population. If we select two samples at random from the same population, we will almost certainly choose different individuals. So the sample results will differ somewhat, just by chance. Properly designed samples avoid systematic bias. But their results are rarely exactly correct, and we expect them to vary from sample to sample.

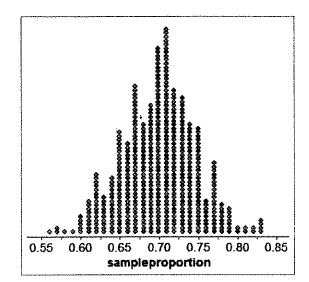
Going to class

How much do sample results vary?

Suppose that 70% of the students in a large university attended all their classes last week. Imagine taking a simple random sample of 100 students and recording the proportion of students in the sample who went to every class last week. Would the sample proportion be exactly 0.70? Probably not. Would the sample proportion be close to 0.70? That depends on what we mean by "close." The following graph shows the results of taking 500 SRSs, each of size 100, and recording the proportion of students who attended all their classes in each sample.

What do we see? The graph is centered at about 0.70, the population proportion. All of the sample proportions fall between 0.55 and 0.85. So we shouldn't be surprised if the difference between the sample proportion and the population proportion is as large as 0.15. The graph also has a very distinctive "bell shape."

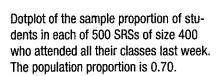


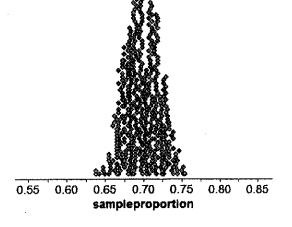


Dotplot of the sample proportion of students in each of 500 SRSs of size 100 who attended all their classes last week. The population proportion is 0.70.

Why can we trust random samples? As the previous example illustrates, the results of random sampling don't change haphazardly from sample to sample. Because we deliberately use chance, the results obey the laws of probability that govern chance behavior. These laws allow us to say how likely it is that sample results are close to the truth about the population. The second reason to use random sampling is that the laws of probability allow trustworthy inference about the population. Results from random samples come with a "margin of error" that sets bounds on the size of the likely error. We will discuss the details of inference for sampling later.

One point is worth making now: larger random samples give better information about the population than smaller samples. For instance, let's look at what happens if we increase the sample size in the example from 100 to 400 students. The dotplot below shows the results of taking 500 SRSs, each of size 400, and recording the proportion of students who attended all their classes in each sample. This graph is also centered at about 0.70. But now all the sample proportions fall between 0.63 and 0.77. So the difference between the sample proportion and the population proportion is at most 0.07. When using SRSs of size 100, this difference could be as much as 0.15. The moral of the story: by taking a very large random sample, you can be confident that the sample result is very close to the truth about the population.





The Current Population Survey contacts about 60,000 households, so we'd expect its estimate of the national unemployment rate to be within about 0.1% of the actual population value. Opinion polls that contact 1000 or 1500 people give less precise results—we expect the sample result to be within about 3% of the actual population percent with a given opinion. Of course, only samples chosen by chance carry this guarantee. Lou Dobbs's online sample tells us little about overall American public opinion even though 7350 people clicked a response.

Sample Surveys: What Can Go Wrong?

The use of bad sampling methods (convenience or voluntary response) often leads to bias. Researchers can avoid bad methods by using random sampling to choose their samples. Other problems in conducting sample surveys are more difficult to avoid.

Sampling is often done using a list of individuals in the population. Such lists are seldom accurate or complete. The result is **undercoverage**.

The list of individuals from which a sample will be drawn is called the sampling frame.

DEFINITION: Undercoverage

Undercoverage occurs when some members of the population cannot be chosen in a sample.

Most samples suffer from some degree of undercoverage. A sample survey of households, for example, will miss not only homeless people but also prison inmates and students in dormitories. An opinion poll conducted by calling landline telephone numbers will miss households that have only cell phones as well as households without a phone. The results of national sample surveys therefore have some bias due to undercoverage if the people not covered differ from the rest of the population.

Well-designed sample surveys avoid bias in the sampling process. The real problems start after the sample is chosen.

One of the most serious sources of bias in sample surveys is nonresponse.

DEFINITION: Nonresponse

Nonresponse occurs when an individual chosen for the sample can't be contacted or refuses to participate.

Nonresponse to surveys often exceeds 50%, even with careful planning and several follow-up calls. If the people who respond differ from those who don't, in a way that is related to the response, bias results.

Some students misuse the term "voluntary response" to explain why certain individuals don't respond in a sample survey. Their idea is that participation in the survey is optional (voluntary), so anyone can refuse to take part. What the students are describing is nonresponse. Think about it this way: nonresponse can occur only after a sample has been selected. In a voluntary response sample, every individual has opted to take part, so there won't be any nonresponse.



EXAMPLE

The ACS, GSS, and Opinion Polls

How bad is nonresponse?

The Census Bureau's American Community Survey (ACS) has the lowest nonresponse rate of any poll we know: only about 1% of the households in the sample refuse to respond. The overall nonresponse rate, including "never at home" and other causes, is just 2.5%. This monthly survey of about 250,000 households replaces the "long form" that in the past was sent to some households in the every-ten-years national census. Participation in the ACS is mandatory, and the Census Bureau follows up by telephone and then in person if a household doesn't return the mail questionnaire.

The University of Chicago's General Social Survey (GSS) is the nation's most important social science survey (see Figure 4.1). The GSS contacts its sample in person, and it is run by a university. Despite these advantages, its most recent survey had a 30% rate of nonresponse.

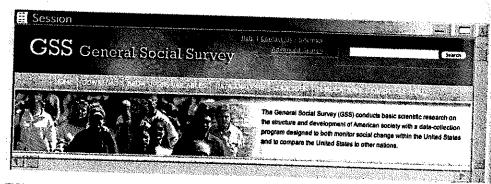
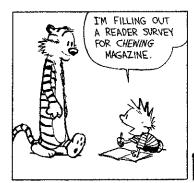


FIGURE 4.1 The home page of the General Social Survey at the University of Chicago's National Opinion Research Center (http://www3.norc.org/GSS+Website/). The GSS has tracked opinions about a wide variety of issues since 1972.

What about opinion polls by news media and opinion-polling firms? We don't know their rates of nonresponse because they won't say. That's a bad sign. The Pew Research Center for the People and the Press imitated a careful random digit dialing survey and published the results: over 5 days, the survey reached 76% of the households in its chosen sample, but "because of busy schedules, skepticism and outright refusals, interviews were completed in just 38% of households that were reached." Combining households that could not be contacted with those who did not complete the interview gave a nonresponse rate of 73%. 11

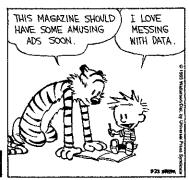
Another type of nonsampling problem occurs when people give inaccurate answers to survey questions. People may lie about their age, income, or drug use. They may misremember how many hours they spent on the Internet last week. Or they might make up an answer to a question that they don't understand.

The gender, race, age, ethnicity, or behavior of the interviewer can also affect people's responses. A systematic pattern of inaccurate answers in a survey leads to response bias.



SEE, THEY ASKED HOW MUCH MONEY I SPEND ON GUM EACH WEEK, SO I WROTE, \$500. FOR MY AGE, I PUT "43, AND WHEN THEY ASKED WHAT MY FAVORITE FLAVOR IS, I WROTE "GARLIC/CURRY."





The wording of questions is the most important influence on the answers given to a sample survey. Confusing or leading questions can introduce strong bias. Changes in wording can greatly affect a survey's outcome.



How Do Americans Feel about Illegal Immigrants?

Question wording matters

"Should illegal immigrants be prosecuted and deported for being in the U.S. illegally, or shouldn't they?" Asked this question in an opinion poll, 69% favored deportation. But when the very same sample was asked whether illegal immigrants who have worked in the United States for two years "should be given a chance to keep their jobs and eventually apply for legal status," 62% said that they should. Different questions give quite different impressions of attitudes toward illegal immigrants.

Even the order in which questions are asked matters. Don't trust the results of a sample survey until you have read the exact questions asked.





Does the order matter? Ask a sample of college students these two questions:

"How happy are you with your life in general?" (Answers on a scale of 1 to 5) "How many dates did you have last month?"

There is almost no association between responses to the two questions when asked in this order. It appears that dating has little to do with happiness. Reverse the order of the questions, however, and a much stronger association appears: college students who say they had more dates tend to give higher ratings of happiness about life. Asking a question that brings dating to mind makes dating success a big factor in happiness.



CHECK YOUR UNDERSTANDING

- 1. Each of the following is a possible source of bias in a sample survey. Name the type of bias that could result.
- (a) The sample is chosen at random from a telephone directory.
- (b) Some people cannot be contacted in five calls.
- (c) Interviewers choose people walking by on the sidewalk to interview.
- 2. A survey paid for by makers of disposable diapers found that 84% of the sample opposed banning disposable diapers. Here is the actual question:

It is estimated that disposable diapers account for less than 2% of the trash in today's landfills. In contrast, beverage containers, third-class mail, and yard wastes are estimated to account for about 21% of the trash in landfills. Given this, in your opinion, would it be fair to ban disposable diapers?¹²

Explain how the wording of the question could result in bias. Be sure to specify the direction of the bias.

Section 4.1 Summary

- A census collects data from every individual in the population.
- A sample survey selects a sample from the population of all individuals about
 which we desire information. The goal of a sample survey is inference: we
 draw conclusions about the population based on data from the sample. It is
 important to specify exactly what population you are interested in and what
 variables you will measure.
- Convenience samples choose individuals who are easiest to reach. In voluntary response samples, individuals choose to join the sample in response to an open invitation. Both of these sampling methods usually lead to bias they consistently underestimate or consistently overestimate the value you want to know.
- Random sampling uses chance to select a sample.
- The basic random sampling method is a simple random sample (SRS). An SRS gives every possible sample of a given size the same chance to be chosen. Choose an SRS by labeling the members of the population and using slips of paper, random digits, or technology to select the sample.
- To choose a stratified random sample, divide the population into stratagroups of individuals that are similar in some way that might affect their responses. Then choose a separate SRS from each stratum and combine these SRSs to form the sample. When strata are "similar within but different between," stratified random samples tend to give more precise estimates of unknown population values than simple random samples.
- To choose a cluster sample, divide the population into groups of individuals that are located near each other, called clusters. Randomly select some of these clusters. All the individuals in the chosen clusters are included in the sample. Ideally, clusters are "different within but similar between." Cluster

- sampling saves time and money by collecting data from entire groups of individuals that are close together.
- Random sampling helps avoid bias in choosing a sample. Bias can still occur
 in the sampling process due to undercoverage, which happens when some
 members of the population cannot be chosen.
- The most serious errors in sample surveys, however, are ones that occur after the sample is chosen. The single biggest problem is **nonresponse**: when people can't be contacted or refuse to answer. Incorrect answers by respondents can lead to **response bias**. Finally, the **wording of questions** has a big influence on the answers.

4.1 TECHNOLOGY CORNER

TI-Nspire instructions in Appendix B; HP Prime instructions on the book's Web site.

10. Choosing an SRS

page 215

Section 4.1 Exercises

- Students as customers A high school's student newspaper plans to survey local businesses about the importance of students as customers. From an alphabetical list of all local businesses, the newspaper staff chooses 150 businesses at random. Of these, 73 return the questionnaire mailed by the staff. Identify the population and the sample.
- Student archaeologists An archaeological dig turns up large numbers of pottery shards, broken stone tools, and other artifacts. Students working on the project classify each artifact and assign it a number. The counts in different categories are important for understanding the site, so the project director chooses 2% of the artifacts at random and checks the students' work. Identify the population and the sample.
- Sampling stuffed envelopes A large retailer prepares its customers' monthly credit card bills using an automatic machine that folds the bills, stuffs them into envelopes, and seals the envelopes for mailing. Are the envelopes completely sealed? Inspectors choose 40 envelopes at random from the 1000 stuffed each hour for visual inspection. Identify the population and the sample.
- Customer satisfaction A department store mails a customer satisfaction survey to people who make credit card purchases at the store. This month,

- 45,000 people made credit card purchases. Surveys are mailed to 1000 of these people, chosen at random, and 137 people return the survey form. Identify the population and the sample.
- 5. Call the shots An advertisement for an upcoming TV show asked: "Should handgun control be tougher? You call the shots in a special call-in poll tonight. If yes, call 1-900-720-6181. If no, call 1-900-720-6182. Charge is 50 cents for the first minute." Over 90% of people who called in said "Yes." Explain why this opinion poll is almost certainly biased.
- 6. Explain it to the congresswoman You are on the staff of a member of Congress who is considering a bill that would provide government-sponsored insurance for nursing-home care. You report that 1128 letters have been received on the issue, of which 871 oppose the legislation. "I'm surprised that most of my constituents oppose the bill. I thought it would be quite popular," says the congresswoman. Are you convinced that a majority of the voters oppose the bill? How would you explain the statistical issue to the congresswoman?
- 7. Instant opinion A recent online poll posed the question "Should female athletes be paid the same as men for the work they do?" In all, 13,147 (44%) said "Yes," 15,182 (51%) said "No," and the

- remaining 1448 said "Don't know." In spite of the large sample size for this survey, we can't trust the result. Why not?
- 8. Sampling at the mall You have probably seen the mall interviewer, approaching people passing by with clipboard in hand. Explain why even a large sample of mall shoppers would not provide a trustworthy estimate of the current unemployment rate.
- 9. Sleepless nights How much sleep do high school students get on a typical school night? An interested student designed a survey to find out. To make data collection easier, the student surveyed the first 100 students to arrive at school on a particular morning. These students reported an average of 7.2 hours of sleep on the previous night.
 - (a) What type of sample did the student obtain?
 - (b) Explain why this sampling method is biased. Is 7.2 hours probably higher or lower than the true average amount of sleep last night for all students at the school? Why?
 - 10. Online polls In June 2008, Parade magazine posed the following question: "Should drivers be banned from using all cell phones?" Readers were encouraged to vote online at parade.com. The July 13, 2008, issue of Parade reported the results: 2407 (85%) said "Yes" and 410 (15%) said "No."
 - (a) What type of sample did the Parade survey obtain?
 - (b) Explain why this sampling method is biased. Is 85% probably higher or lower than the true percent of all adults who believe that cell phone use while driving should be banned? Why?
- 11. Do you trust the Internet? You want to ask a sample of high school students the question "How much do you trust information about health that you find on the Internet—a great deal, somewhat, not much, or not at all?" You try out this and other questions on a pilot group of 5 students chosen from your class. The class members are listed below.
 - (a) Explain how you would use a line of Table D to choose an SRS of 5 students from the following list. Explain your method clearly enough for a classmate to obtain your results.
 - (b) Use line 107 to select the sample. Show how you use each of the digits.

| | | · · · · · · · · · · · · · · · · · · · | | |
|----------|-----------|---------------------------------------|-----------|----------|
| Anderson | Deng | Glaus | Nguyen | Samuels |
| Arroyo | De Ramos | Helling | Palmiero | Shen |
| Batista | Drasin | Husain | Percival | Tse |
| Bell | Eckstein | Johnson | Prince | Velasco |
| Burke | Fernandez | Kim | Puri | Wallace |
| Cabrera | Fullmer | Molina | Richards | Washburn |
| Calloway | Gandhi | Morgan | Rider | Zabidi |
| Delluci | Garcia | Murphy | Rodriguez | Zhao |
| | | | | |

- 12. Apartment living You are planning a report on apartment living in a college town. You decide to select three apartment complexes at random for indepth interviews with residents.
- (a) Explain how you would use a line of Table D to choose an SRS of 3 complexes from the list below. Explain your method clearly enough for a classmate to obtain your results.
- (b) Use line 117 to select the sample. Show how you use each of the digits.

| Ashley Oaks | Chauncey Village | Franklin Park | Richfield |
|---------------|------------------|-----------------|------------------|
| Bay Pointe | Country Squire | Georgetown | Sagamore Ridge |
| Beau Jardin | Country View | Greenacres | Salem Courthouse |
| Bluffs | Country Vilta | Lahr House | Village Manor |
| Brandon Place | Crestview | Mayfair Village | Waterford Court |
| Briarwood | Del-Lynn | Nobb Hill | Williamsburg |
| Brownstone | Fairington | Pemberly Courts | - ; |
| Burberry | Fairway Knolls | Peppermill | |
| Cambridge | Fowler | Pheasant Run | : |

- 13. Sampling the forest To gather data on a 1200-acre pine forest in Louisiana, the U.S. Forest Service laid a grid of 1410 equally spaced circular plots over a map of the forest. A ground survey visited a sample of 10% of these plots. 13
- (a) Explain how you would use your calculator or Table D to choose an SRS of 141 plots. Your description should be clear enough for a classmate to carry out your plan.
- (b) Use your method from (a) to choose the first 3 plots.
- 14. Sampling gravestones The local genealogical society in Coles County, Illinois, has compiled records on all 55,914 gravestones in cemeteries in the county for the years 1825 to 1985. Historians plan to use these records to learn about African Americans in Coles County's history. They first choose an SRS of 395 records to check their accuracy by visiting the actual gravestones. 14
- (a) Explain how you would use your calculator or Table D to choose the SRS. Your description should be clear enough for a classmate to carry out your plan.
- (b) Use your method from (a) to choose the first 3 gravestones.
- 15. Random digits Which of the following statements are true of a table of random digits, and which are false? Briefly explain your answers.
- (a) There are exactly four 0s in each row of 40 digits.
- (b) Each pair of digits has chance 1/100 of being 00.
- (c) The digits 0000 can never appear as a group, because this pattern is not random.
- 16. Random digits In using Table D repeatedly to choose random samples, you should not always begin at the same place, such as line 101. Why not?

- iPhones Suppose 1000 iPhones are produced at a factory today. Management would like to ensure that the phones' display screens meet their quality standards before shipping them to retail stores. Since it takes about 10 minutes to inspect an individual phone's display screen, managers decide to inspect a sample of 20 phones from the day's production.
- Explain why it would be difficult for managers to inspect an SRS of 20 iPhones that are produced today.
- An eager employee suggests that it would be easy to inspect the last 20 iPhones that were produced today. Why isn't this a good idea?
- Another employee recommends a different sampling method: Randomly choose one of the first 50 iPhones produced. Inspect that phone and every fiftieth iPhone produced afterward. (This method is known as systematic random sampling.) Explain carefully why this sampling method is not an SRS.
- 18. Dead trees On the west side of Rocky Mountain National Park, many mature pine trees are dying due to infestation by pine beetles. Scientists would like to use sampling to estimate the proportion of all pine trees in the area that have been infested.
- Explain why it wouldn't be practical for scientists to obtain an SRS in this setting.
- A possible alternative would be to use every pine tree along the park's main road as a sample. Why is this sampling method biased?
- Suppose that a more complicated random sampling plan is carried out, and that 35% of the pine trees in the sample are infested by the pine beetle. Can scientists conclude that exactly 35% of all the pine trees on the west side of the park are infested? Why or why not?
- Who goes to the convention? A club has 30 student members and 10 faculty members. The students are

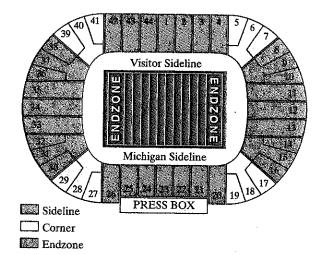
| Abel | Fisher | Huber | Miranda | Reinmann |
|----------|-----------|---------|-----------|----------|
| Carson | Ghosh | Jimenez | Moskowitz | Santos |
| Chen | Griswold | Jones | Neyman | Shaw |
| David | Hein | Kim | O'Brien | Thompson |
| Deming | Hernandez | Klotz | Pearl | Utts |
| Elashoff | Holland | Liu | Potter | Varga |

The faculty members are

| Andrews | Fernandez | Kim | Moore | West |
|-------------|-----------|----------|----------|------|
| Besicovitch | Gupta | Lightman | Phillips | Yang |

The club can send 4 students and 2 faculty members to a convention. It decides to choose those who will go by random selection. Describe a method for using Table D to select a stratified random sample of 4 students and 2 faculty. Then use line 123 to select the sample.

- 20. Sampling by accountants Accountants often use stratified samples during audits to verify a company's records of such things as accounts receivable. The stratification is based on the dollar amount of the item and often includes 100% sampling of the largest items. One company reports 5000 accounts receivable. Of these, 100 are in amounts over \$50,000; 500 are in amounts between \$1000 and \$50,000; and the remaining 4400 are in amounts under \$1000. Using these groups as strata, you decide to verify all the largest accounts and to sample 5% of the midsize accounts and 1% of the small accounts. Describe a method for using Table D to select a stratified random sample of the midsize and small accounts. Then use line 115 to select only the first 3 accounts from each of these strata.
 - Go Blue! Michigan Stadium, also known as "The Big House," seats over 100,000 fans for a football game. The University of Michigan athletic department plans to conduct a survey about concessions that are sold during games. Tickets are most expensive for seats on the sidelines. The cheapest seats are in the end zones (where one of the authors sat as a student). A map of the stadium is shown.



- (a) The athletic department is considering a stratified random sample. What would you recommend as the strata? Why?
- Explain why a cluster sample might be easier to obtain. What would you recommend for the clusters? Why?
- 22. How was your stay? A hotel has 30 floors with 40 rooms per floor. The rooms on one side of the hotel face the water, while rooms on the other side face a golf course. There is an extra charge for the rooms with a water view. The hotel manager wants to survey 120 guests who stayed at the hotel during a convention about their overall satisfaction with the property.

(a) Explain why choosing a stratified random sample might be preferable to an SRS in this case. What would you use as strata?

232

- (b) Why might a cluster sample be a simpler option? What would you use as clusters?
- 23. Is it an SRS? A corporation employs 2000 male and 500 female engineers. A stratified random sample of 200 male and 50 female engineers gives each engineer 1 chance in 10 to be chosen. This sample design gives every individual in the population the same chance to be chosen for the sample. Is it an SRS? Explain your answer.
- 24. Attitudes toward alcohol At a party there are 30 students over age 21 and 20 students under age 21. You choose at random 3 of those over 21 and separately choose at random 2 of those under 21 to interview about attitudes toward alcohol. You have given every student at the party the same chance to be interviewed. Why is your sample not an SRS?
- 25. High-speed Internet Laying fiber-optic cable is expensive. Cable companies want to make sure that if they extend their lines out to less dense suburban or rural areas, there will be sufficient demand and the work will be cost-effective. They decide to conduct a survey to determine the proportion of households in a rural subdivision that would buy the service. They select a simple random sample of 5 blocks in the subdivision and survey each family that lives on one of those blocks.
- (a) What is the name for this kind of sampling method?
- (b) Give a possible reason why the cable company chose this method.
- 26. Timber! A lumber company wants to estimate the proportion of trees in a large forest that are ready to be cut down. They use an aerial map to divide the forest into 200 equal-sized rectangles. Then they choose a random sample of 20 rectangles and examine every tree that's in one of those rectangles.
- (a) What is the name for this kind of sampling method?
- (b) Give a possible reason why the lumber company chose this method.
- 27. Tweet, tweet! What proportion of students at your school use Twitter? To find out, you survey a simple random sample of students from the school roster.
- (a) Will your sample result be exactly the same as the true population proportion? Explain.
- (b) Which would be more likely to get your sample result closer to the true population value: an SRS of 50 students or an SRS of 100 students? Explain.
- 28. Far from home? A researcher wants to estimate the average distance that students at a large community college live from campus. To find out, she surveys a simple random sample of students from the registrar's database.

- (a) Will the researcher's sample result be exactly the same as the true population mean? Explain.
- (b) Which would be more likely to get the researcher's sample result closer to the true population value: an SRS of 100 students or an SRS of 50 students? Explain.
- 29. Baseball tickets Suppose you want to know the average amount of money spent by the fans attending opening day for the Cleveland Indians baseball season. You get permission from the team's management to conduct a survey at the stadium, but they will not allow you to bother the fans in the club seating or box seats (the most expensive seating). Using a computer, you randomly select 500 seats from the rest of the stadium. During the game, you ask the fans in those seats how much they spent that day.

Give a reason why this survey might yield a biased result. Explain the likely direction of the bias.

30. Rise and shine How long before school starts do students get out of bed, on average? Administrators survey a random sample of students on each school bus one morning.

Give a reason why this survey might yield a biased result. Explain the likely direction of the bias.

- 31. Nonresponse A survey of drivers began by randomly sampling all listed residential telephone numbers in the United States. Of 45,956 calls to these numbers, 5029 were completed. The goal of the survey was to estimate how far people drive, on average, per day.¹⁵
- (a) What was the rate of nonresponse for this sample?
- (b) Explain how nonresponse can lead to bias in this survey. Be sure to give the direction of the bias.
- 32. Ring-no-answer A common form of nonresponse in telephone surveys is "ring-no-answer." That is, a call is made to an active number but no one answers. The Italian National Statistical Institute looked at nonresponse to a government survey of households in Italy during the periods January 1 to Easter and July 1 to August 31. All calls were made between 7 and 10 P.M., but 21.4% gave "ring-no-answer" in one period versus 41.5% "ring-no-answer" in the other period. Which period do you think had the higher rate of no answers? Why? Explain why a high rate of nonresponse makes sample results less reliable.
- 33. Running red lights The sample described in Exercise 31 produced a list of 5024 licensed drivers. The investigators then chose an SRS of 880 of these drivers to answer questions about their driving habits. One question asked was: "Recalling the last ten traffic lights you drove through, how many of them were red when you entered the intersections?" Of the 880 respondents, 171 admitted that at least one light had been red. A practical problem with this survey is that

Section 4.1 Sampling and Surveys



- people may not give truthful answers. What is the likely direction of the bias? Explain.
- 34. Seat belt use A study in El Paso, Texas, looked at seat belt use by drivers. Drivers were observed at randomly chosen convenience stores. After they left their cars, they were invited to answer questions that included questions about seat belt use. In all, 75% said they always used seat belts, yet only 61.5% were wearing seat belts when they pulled into the store parking lots. The Explain the reason for the bias observed in responses to the survey. Do you expect bias in the same direction in most surveys about seat belt use?
- 35. Wording bias Comment on each of the following as a potential sample survey question. Is the question clear? Is it slanted toward a desired response?
- (a) "Some cell phone users have developed brain cancer. Should all cell phones come with a warning label explaining the danger of using cell phones?"
- (b) "Do you agree that a national system of health insurance should be favored because it would provide health insurance for everyone and would reduce administrative costs?"
- (c) "In view of escalating environmental degradation and incipient resource depletion, would you favor economic incentives for recycling of resource intensive consumer goods?"
- 36. Checking for bias Comment on each of the following as a potential sample survey question. Is the question clear? Is it slanted toward a desired response?
- (a) Which of the following best represents your opinion on gun control?
 - 1. The government should confiscate our guns.
 - 2. We have the right to keep and bear arms.
- (b) A freeze in nuclear weapons should be favored because it would begin a much-needed process to stop everyone in the world from building nuclear weapons now and reduce the possibility of nuclear war in the future. Do you agree or disagree?

Multiple choice: Select the best answer for Exercises 37 to 42.

- 37. The Web portal AOL places opinion poll questions next to many of its news stories. Simply click your response to join the sample. One of the questions in January 2008 was "Do you plan to diet this year?" More than 30,000 people responded, with 68% saying "Yes." You can conclude that
- (a) about 68% of Americans planned to diet in 2008.
- (b) the poll used a convenience sample, so the results tell us little about the population of all adults.
- the poll uses voluntary response, so the results tell us little about the population of all adults.

- (d) the sample is too small to draw any conclusion.
- (e) None of these.
- 38. To gather information about the validity of a new standardized test for tenth-grade students in a particular state, a random sample of 15 high schools was selected from the state. The new test was administered to every 10th-grade student in the selected high schools. What kind of sample is this?
- (a) A simple random sample
- (b) A stratified random sample
- (c) A cluster sample
- (d) A systematic random sample
- (e) A voluntary response sample
- 39. Your statistics class has 30 students. You want to call an SRS of 5 students from your class to ask where they use a computer for the online quizzes. You label the students 01, 02, ..., 30. You enter the table of random digits at this line:

14459 26056 31424 80371 65103 62253 22490 61181 Your SRS contains the students labeled

- (a) 14, 45, 92, 60, 56.
- **(b)** 14, 31, 03, 10, 22.
- (c) 14, 03, 10, 22, 22.
- (d) 14, 03, 10, 22, 06.
- (e) 14, 03, 10, 22, 11.
- 40. Suppose that 35% of the registered voters in a state are registered as Republicans, 40% as Democrats, and 25% as Independents. A newspaper wants to select a sample of 1000 registered voters to predict the outcome of the next election. If they randomly select 350 Republicans, randomly select 400 Democrats, and randomly select 250 Independents, did this sampling procedure result in a simple random sample of registered voters from this state?
- (a) Yes, because each registered voter had the same chance of being chosen.
- (b) Yes, because random chance was involved.
- (c) No, because not all registered voters had the same chance of being chosen.
- (d) No, because there were a different number of registered voters selected from each party.
- (e) No, because not all possible groups of 1000 registered voters had the same chance of being chosen.
- 41. A local news agency conducted a survey about unemployment by randomly dialing phone numbers until they had gathered responses from 1000 adults in their state. In the survey, 19% of those who responded said they were not currently employed. In reality, only 6% of the adults in the state were not currently employed

- at the time of the survey. Which of the following best explains the difference in the two percentages?
- (a) The difference is due to sampling variability. We shouldn't expect the results of a random sample to match the truth about the population every time.

234

- (b) The difference is due to response bias. Adults who are employed are likely to lie and say that they are unemployed.
- (c) The difference is due to undercoverage bias. The survey included only adults and did not include teenagers who are eligible to work.
- (d) The difference is due to nonresponse bias. Adults who are employed are less likely to be available for the sample than adults who are unemployed.
- (e) The difference is due to voluntary response. Adults are able to volunteer as a member of the sample.
- 42. A simple random sample of 1200 adult Americans is selected, and each person is asked the following question: "In light of the huge national deficit, should the government at this time spend additional money to establish a national system of health insurance?" Only 39% of those responding answered "Yes." This survey
- (a) is reasonably accurate since it used a large simple random sample.
- (b) needs to be larger since only about 24 people were drawn from each state.
- (c) probably understates the percent of people who favor a system of national health insurance.

- (d) is very inaccurate but neither understates nor overstates the percent of people who favor a system of national health insurance. Because simple random sampling was used, it is unbiased.
- (e) probably overstates the percent of people who favor a system of national health insurance.
- 43. Sleep debt (3.2) A researcher reported that the typical teenager needs 9.3 hours of sleep per night but gets only 6.3 hours. ¹⁸ By the end of a 5-day school week, a teenager would accumulate about 15 hours of "sleep debt." Students in a high school statistics class were skeptical, so they gathered data on the amount of sleep debt (in hours) accumulated over time (in days) by a random sample of 25 high school students. The resulting least-squares regression equation for their data is Sleep debt = 2.23 + 3.17(days).
- (a) Interpret the slope of the regression line in context.
- (b) Are the students' results consistent with the researcher's report? Explain.
- 44. Internet charges (2.1) Some Internet service providers (ISPs) charge companies based on how much bandwidth they use in a month. One method that ISPs use for calculating bandwidth is to find the 95th percentile of a company's usage based on samples of hundreds of 5-minute intervals during a month.
- (a) Explain what "95th percentile" means in this setting.
- (b) Which would cost a company more: the 95th percentile method or a similar approach using the 98th percentile? Justify your answer.

4.2 Experiments

WHAT YOU WILL LEARN

By the end of the section, you should be able to:

- Distinguish between an observational study and an experiment.
- Explain the concept of confounding and how it limits the ability to make cause-and-effect conclusions.
- Identify the experimental units, explanatory and response variables, and treatments in an experiment.
- Explain the purpose of comparison, random assignment, control, and replication in an experiment.

- Describe a completely randomized design for an experiment, including how to randomly assign treatments using slips of paper, technology, or a table of random digits.
- Describe the placebo effect and the purpose of blinding in an experiment.
- Interpret the meaning of statistically significant in the context of an experiment.
- Explain the purpose of blocking in an experiment.
 Describe a randomized block design or a matched pairs design for an experiment.

R3.5 (a) $\hat{y} = 30.2 + 0.16x$, where y = final exam score and x = total score before the final examination. (b) 78.2 (c) Of all the lines that the professor could use to summarize the relationship between final exam score and total points before the final exam, the least-squares regression line is the one that has the smallest sum of squared residuals. (d) Because $r^2 = 0.36$, only 36% of the variability in the final exam scores is accounted for by the linear model relating final exam scores to total score before the final exam. More than half (64%) of the variation in final exam scores is not accounted for, so Julie has reason to question this estimate.

R3.6 Even though there is a high correlation between number of calculators and math achievement, we shouldn't conclude that increasing the number of calculators will *cause* an increase in math achievement. It is possible that students who are more serious about school have better math achievement and also have more calculators.

Answers to Chapter 3 AP® Statistics Practice Test

T3.1 d

Т3.2 е

Т3.3 с

T3.4 a

T3.5 a

Т3.6 с

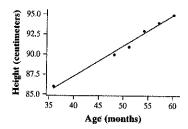
Т3.7 Ъ

T3.8 e

Т3.9 Ъ

T3.10 c

T3.11 (a) A scatterplot with regression line is shown below. (b) $\hat{y} = 71.95 + 0.3833x$, where y = height and x = age. (c) 255.934 cm, or 100.76 inches (d) This was an extrapolation. Our data were based only on the first 5 years of life and the linear trend will not continue forever.



T3.12 (a) The point in the upper-right-hand corner has a very high silicon value for its isotope value. (b) (i) r would get closer to -1 because it does not follow the linear pattern of the other points. (ii) Because this point is "pulling up" the line on the right side of the plot, removing it will make the slope steeper (more negative) and the y intercept smaller (note that the y axis is to the right of the points in the scatterplot). (iii) Because this point has a large residual, removing it will make s a little smaller.

T3.13 (a) $\hat{y} = 92.29 - 0.05762x$, where y is the percent of the grass burned and x is the number of wildebeest. (b) For every increase of 1000 wildebeest, the predicted percent of grassy area burned decreases by about 0.058. (c) $r = -\sqrt{0.646} = -0.804$. There is a strong, negative linear association between the percent of grass burned and the number of wildebeest. (d) Yes, because there is no obvious leftover pattern in the residual plot.

Chapter 4

Section 4.1

Answers to Check Your Understanding

page 213: I. Convenience sampling. This could lead the inspector to overestimate the quality of the oranges if the farmer puts the best oranges on top. 2. Voluntary response sampling. In this case, those who are happy that the UN has its headquarters in the U.S. already have what they want and so are less likely to respond. The proportion who answered "No" in the sample is likely to be higher than the true proportion in the U.S. who would answer "No."

page 223: 1. You would have to identify 200 different seats, go to those seats in the arena, and find the people who are sitting there, which would take a lot of time. 2. It is best to create strata where the people within a stratum are very similar to each other but different than the people in other strata. In this case, it would be better to take the lettered rows as the strata because each lettered row is the same distance from the court and so would contain only seats with the same (or nearly the same) ticket price. 3. It is best if the people in each cluster reflect the variability found in the population. In this case, it would be better to take the numbered sections as the clusters because they include all different seat prices.

page 228: 1. (a) Undercoverage (b) Nonresponse (c) Undercoverage 2. By making it sound like they are not a problem in the landfill, this question will result in fewer people suggesting that we should ban disposable diapers. The proportion who would say "Yes" to this survey question is likely to be smaller than the proportion who would say "Yes" to a more fairly worded question.

Answers to Odd-Numbered Section 4.1 Exercises

4.1 Population: all local businesses. Sample: the 73 businesses that return the questionnaire.

4.3 Population: the 1000 envelopes stuffed during a given hour. Sample: the 40 randomly selected envelopes.

4.5 This is a voluntary response sample. In this case, it appears that people who strongly support gun control volunteered more often, causing the proportion in the sample to be greater than the proportion in the population.

4.7 This is a voluntary response sample and overrepresents the opinions of those who feel most strongly about the issue being surveyed.

4.9 (a) A convenience sample (b) The first 100 students to arrive at school likely had to wake up earlier than other students, so 7.2 hours is probably less than the true average.

4.11 (a) Number the 40 students from 01 to 40. Pick a starting point on the random number table. Record two-digit numbers, skipping numbers that aren't between 01 and 40 and any repeated numbers, until you have 5 unique numbers between 01 and 40. Use the 5 students corresponding to these numbers. (b) Using line 107, skip the numbers not in bold: 82 73 95 78 90 20 80 74 75 11 81 67 65 53 00 94 38 31 48 93 60 94 07. Select Johnson (20), Drasin (11), Washburn (38), Rider (31), and Calloway (07).

4.13 (a) Using calculator: Number the plots from 1 to 1410. Use the command randInt (1,1410) to select 141 different integers from 1 to 1410 and use the corresponding 141 plots. (b) Answers will vary.

4.15 (a) False—although, on average, there will be four 0s in every set of 40 digits, the number of 0s can be less than 4 or greater than 4 by chance. (b) True—there are 100 pairs of digits 00 through 99,

and all are equally likely. (c) False – 0000 is just as likely as any other string of four digits.

4.17 (a) It might be difficult to locate the 20 phones from among the 1000 produced that day. (b) The quality of the phones produced may change during the day, so that the last phones manufactured are not representative of the day's production. (c) Because each sample of 20 phones does not have the same probability of being selected. In an SRS, it is possible for 2 consecutive phones to be selected in a sample, but this is not possible with a systematic random sample.

4.19 Assign numbers 01 to 30 to the students. Pick a starting point on the random digit table. Record two-digit numbers, skipping any that aren't between 01 and 30 and any repeated numbers, until you have 4 unique numbers between 01 and 30. Use the corresponding four students. Then assign numbers 0 to 9 to the faculty members. Continuing on the table, record one-digit numbers, skipping any repeated numbers, until you have 2 unique numbers between 0 and 9. Use the corresponding faculty members. Starting on line 123 gives 08-Ghosh, 15-Jones, 07-Fisher, and 27-Shaw for the students and 1-Besicovitch and 0-Andrews for the faculty.

4.21 (a) Use the three types of seats as the strata because people who can afford more expensive tickets probably have different opinions about the concessions than people who can afford only the cheaper tickets. (b) A stratified random sample will include seats from all over the stadium, which would make it very time-consuming to obtain. A cluster sample of numbered sections would be easier to obtain, because the people selected for the sample would be sitting close together.

4.23 No. In an SRS, each possible sample of 250 engineers is equally likely to be selected, including samples that aren't exactly 200 males and 50 females.

4.25 (a) Cluster sampling. (b) To save time and money. In an SRS, the company would have to visit individual homes all over the rural subdivision instead of only 5 locations.

4.27 (a) It is unlikely, because different random samples will include different students and produce different estimates of the proportion of students who use Twitter. (b) An SRS of 100 students. Larger random samples give us better information about the population than smaller random samples.

4.29 Because you are sampling only from the lower-priced ticket holders, this will likely produce an estimate that is too small, as fans in the club seats and box seats probably spend more money at the game than fans in cheaper seats.

4.31 (a) 89.1% (b) Because the people who have long commutes are less likely to be at home and be included in the sample, this will likely produce an estimate that is too small.

4.33 We would not expect very many people to claim they have run red lights when they haven't, but some people will deny running red lights when they have. Thus, we expect that the sample proportion underestimates the true proportion of drivers who have run a red light.

4.35 (a) The wording is clear, but the question is slanted in favor of warning labels because of the first sentence stating that some cell phone users have developed brain cancer. (b) The question is clear, but it is slanted in favor of national health insurance by asserting it would reduce administrative costs and not providing any counterarguments. (c) The wording is too technical for many people to understand. For those who do understand the question, it is slanted because it suggests reasons why one should support recycling.

4.39 d

4.41 d

4.43 (a) For each additional day, the predicted sleep debt increases by about 3.17 hours. (b) The predicted sleep debt for a 5-day school week is 2.23 + 3.17(5) = 18.08 hours. This is about 3 hours more than the researcher claimed for a 5-day week, so the students have reason to be skeptical of the research study's reported results.

Section 4.2

Answers to Check Your Understanding

page 237: 1. Experiment, because a treatment (brightness of screen) was imposed on the laptops. 2. Observational study, because students were not assigned to eat a particular number of meals with their family per week. 3. Explanatory: number of meals per week eaten with their family. Response: GPA. 4. There are probably other variables that are influencing the response variable. For example, students who have part-time jobs may not be able to eat many meals with their families and may not have much time to study, leading to lower grades.

page 247: 1. Randomly assign the 29 students to two treatments: evaluating the performance in small groups or evaluating the performance alone. The response variable will be the accuracy of their final performance evaluations. To implement this design, use 29 equally sized slips of paper. Label 15 of them "small group" and 14 of them "alone." Then shuffle the papers and hand them out at random to the 29 students, assigning them to a treatment. 2. The purpose of the control group is to provide a baseline for comparison. Without a group to compare to, it is impossible to determine if the small group treatment is more effective.

page 249: 1. No. Perhaps seeing the image of their unborn child encouraged the mothers who had an ultrasound to eat a better diet, resulting in healthier babies. 2. No. While the people weighing the babies at birth may not have known whether that particular mother had an ultrasound or not, the mothers knew. This might have affected the outcome because the mothers knew whether they had received the treatment or not. 3. Treat all mothers as if they had an ultrasound, but for some mothers the ultrasound machine wouldn't be turned on. To avoid having mothers know the machine was turned off, the ultrasound screen would have to be turned away from all the mothers.

Answers to Odd-Numbered Section 4.2 Exercises

4.45 Experiment, because students were randomly assigned to the different teaching methods.

4.47 (a) Observational study, because mothers weren't assigned to eat different amounts of chocolate. (b) Explanatory: the mother's chocolate consumption. Response: the baby's temperament. (c) No, this study is an observational study so we cannot draw a cause-and-effect conclusion. It is possible that women who eat chocolate daily have less stressful lives and the lack of stress helps their babies to have better temperaments.

4.49 Type of school. For example, private schools tend to have smaller class sizes and students that come from families with higher socioeconomic status. If these students do better in the future, we wouldn't know if the better performance was due to smaller class sizes or higher socioeconomic status.

4.51 Experimental units: pine seedlings. Explanatory variable: light intensity. Response variable: dry weight at the end of the study. Treatments: full light, 25% light, and 5% light.

| | • | |
|---|---|--|
| | | |
| | | |
| | | |
| | | |
| · | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |