

THIRD EDITION

Even You Can Learn
STATISTICS
— and —
ANALYTICS

An Easy to Understand Guide
to Statistics and Analytics

David M. Levine and David F. Stephan

Even You Can Learn Statistics and Analytics

Third Edition

**An Easy to Understand Guide to
Statistics and Analytics**

David M. Levine, Ph.D.

David F. Stephan

Editor-in-Chief: Amy Neidlinger
Operations Specialist: Jodi Kemper
Cover Designer: Alan Clements
Managing Editor: Kristy Hart
Senior Project Editor: Betsy Gratner
Copy Editor: Krista Hansing
Proofreader: Sarah Kearns
Interior Designer: Argosy
Compositor: codeMantra
Manufacturing Buyer: Dan Uhrig

© 2015 by Pearson Education, Inc.
Upper Saddle River, New Jersey 07458

For information about buying this title in bulk quantities, or for special sales opportunities (which may include electronic versions; custom cover designs; and content particular to your business, training goals, marketing focus, or branding interests), please contact our corporate sales department at corpsales@pearsoned.com or (800) 382-3419.

For government sales inquiries, please contact governmentsales@pearsoned.com.

For questions about sales outside the U.S., please contact international@pearsoned.com.

Company and product names mentioned herein are the trademarks or registered trademarks of their respective owners.

All rights reserved. No part of this book may be reproduced, in any form or by any means, without permission in writing from the publisher.

Printed in the United States of America

First Printing December 2014

ISBN-10: 0-13-338266-4

ISBN-13: 978-0-13-338266-2

Pearson Education LTD.
Pearson Education Australia PTY, Limited.
Pearson Education Singapore, Pte. Ltd.
Pearson Education North Asia, Ltd.
Pearson Education Canada, Ltd.
Pearson Educación de Mexico, S.A. de C.V.
Pearson Education—Japan
Pearson Education Malaysia, Pte. Ltd.

Library of Congress Control Number: 2014949421

*To our wives
Marilyn and Mary*

*To our children
Sharyn and Mark*

*And to our parents
In loving memory, Lee, Reuben, Ruth, and Francis*

This page intentionally left blank

Table of Contents

Acknowledgments	viii
About the Authors	ix
Introduction <i>The Even You Can Learn Statistics and Analytics Owners Manual</i>	xi
Chapter 1 Fundamentals of Statistics	1
1.1 The First Three Words of Statistics	2
1.2 The Fourth and Fifth Words	4
1.3 The Branches of Statistics	4
1.4 Sources of Data	5
1.5 Sampling Concepts	7
1.6 Sample Selection Methods	9
Chapter 2 Presenting Data in Tables and Charts	15
2.1 Presenting Categorical Variables	15
2.2 Presenting Numerical Variables	22
2.3 “Bad” Charts	28
Chapter 3 Descriptive Statistics	37
3.1 Measures of Central Tendency	37
3.2 Measures of Position	41
3.3 Measures of Variation	45
3.4 Shape of Distributions	51
Chapter 4 Probability	67
4.1 Events	67
4.2 More Definitions	68
4.3 Some Rules of Probability	70
4.4 Assigning Probabilities	73
Chapter 5 Probability Distributions	79
5.1 Probability Distributions for Discrete Variables	79
5.2 The Binomial and Poisson Probability Distributions	85
5.3 Continuous Probability Distributions and the Normal Distribution	92
5.4 The Normal Probability Plot	100

Chapter 6 Sampling Distributions and Confidence Intervals	113
6.1 Foundational Concepts	113
6.2 Sampling Error and Confidence Intervals	117
6.3 Confidence Interval Estimate for the Mean Using the t Distribution (σ Unknown)	120
6.4 Confidence Interval Estimation for Categorical Variables	123
6.5 Bootstrapping Estimation	126
Chapter 7 Fundamentals of Hypothesis Testing	137
7.1 The Null and Alternative Hypotheses	137
7.2 Hypothesis Testing Issues	139
7.3 Decision-Making Risks	141
7.4 Performing Hypothesis Testing	143
7.5 Types of Hypothesis Tests	144
Chapter 8 Hypothesis Testing: Z and t Tests	149
8.1 Testing for the Difference Between Two Proportions	149
8.2 Testing for the Difference Between the Means of Two Independent Groups	156
8.3 The Paired t Test	162
Chapter 9 Hypothesis Testing: Chi-Square Tests and the One-Way Analysis of Variance (ANOVA)	175
9.1 Chi-Square Test for Two-Way Cross-Classification Tables	175
9.2 One-Way Analysis of Variance (ANOVA): Testing for the Differences Among the Means of More Than Two Groups	182
Chapter 10 Simple Linear Regression	203
10.1 Basics of Regression Analysis	203
10.2 Developing a Simple Linear Regression Model	206
10.3 Measures of Variation	215
10.4 Inferences About the Slope	220
10.5 Common Mistakes Using Regression Analysis	223
Chapter 11 Multiple Regression	239
11.1 The Multiple Regression Model	239
11.2 Coefficient of Multiple Determination	242
11.3 The Overall F Test	243
11.4 Residual Analysis for the Multiple Regression Model	244
11.5 Inferences Concerning the Population Regression Coefficients	245
Chapter 12 Fundamentals of Analytics	257
12.1 Basic Vocabulary of Analytics	257
12.2 Software for Analytics	260

Chapter 13 Descriptive Analytics	265
13.1 Dashboards	265
13.2 Common Descriptive Analytics Visualizations	268
Chapter 14 Predictive Analytics	277
14.1 Analysis with Predictive Analytics	277
14.2 Classification and Regression Trees	278
14.3 Cluster Analysis	283
14.4 Multidimensional Scaling	286
Appendix A Microsoft Excel Operation and Configuration	293
A.S1 Spreadsheet Operation Conventions	293
A.S2 Spreadsheet Technical Configurations	294
Appendix B Review of Arithmetic and Algebra	295
Assessment Quiz	295
Symbols	298
Answers to Quiz	304
Appendix C Statistical Tables	305
Appendix D Spreadsheet Tips	333
CT: Chart Tips	333
FT: Function Tips	335
Appendix E Advanced Techniques	337
ADV: Advanced How-Tos	337
ATT: Analysis ToolPak Tips	342
Appendix F Documentation for Downloadable Files	345
F1 Downloadable Data Files	345
F2 Downloadable <i>Spreadsheet Solution</i> Files	348
Glossary	349
Index	357

Acknowledgments

We would especially like to thank the staff at Financial Times/Pearson: Amy Neidlinger for making this book a reality, Sarah Kearns for her proofreading, Krista Hansing for her copy editing, and Betsy Gratner for her work in the production of this text.

We have sought to make the contents of this book as clear, accurate, and error-free as possible. We invite you to make suggestions or ask questions about the content if you think we have fallen short of our goals in any way. Please email your comments to davidlevine@davidlevinestatics.com and include “Even You Can Learn Statistics and Analytics 3/e” in the subject line.

About the Authors

David M. Levine is Professor Emeritus of Statistics and Computer Information Systems at Baruch College-CUNY. He received B.B.A. and M.B.A. degrees in Statistics from City College of New York and a Ph.D. degree from New York University in Industrial Engineering and Operations Research. He is nationally recognized as a leading innovator in business statistics education and is the coauthor of such best-selling statistics textbooks as *Statistics for Managers Using Microsoft Excel*, *Basic Business Statistics: Concepts and Applications*, *Business Statistics: A First Course*, and *Applied Statistics for Engineers and Scientists Using Microsoft Excel and Minitab*.

He also is the author of *Statistics for Six Sigma Green Belts and Champions*, published by Financial Times–Prentice-Hall. He is coauthor of *Six Sigma for Green Belts and Champions* and *Design for Six Sigma for Green Belts and Champions* also published by Financial Times–Prentice-Hall, and *Quality Management* Third Ed., McGraw-Hill-Irwin. He is also the author of *Video Review of Statistics* and *Video Review of Probability*, both published by Video Aided Instruction. He has published articles in various journals including *Psychometrika*, *The American Statistician*, *Communications in Statistics*, *Multivariate Behavioral Research*, *Journal of Systems Management*, *Quality Progress*, and *The American Anthropologist* and has given numerous talks at American Statistical Association, Decision Sciences Institute, and Making Statistics More Effective in Schools of Business conferences. While at Baruch College, Dr. Levine received numerous awards for outstanding teaching.

David F. Stephan is an independent instructional technologist. During his more than 20 years teaching at Baruch College-CUNY, he pioneered the use of computer-equipped classrooms and interdisciplinary multimedia tools, was an associate director of a U.S. Department of Education FIPSE project that applied interactive media to support instruction and devised techniques for teaching computer applications in a business context. A frequent participant in the Decision Sciences Institute's Making Statistics for More Effective in School of Business mini-conferences, he is also a coauthor of *Business Statistics: A First Course* and *Statistics for Managers Using Microsoft Excel*. He is also the developer of PHStat, the statistics add-in for Microsoft Excel distributed by Pearson Education.

This page intentionally left blank

Introduction

The Even You Can Learn Statistics and Analytics Owners Manual

In today's world, understanding statistics and analytics is more important than ever. *Even You Can Learn Statistics and Analytics: A Guide for Everyone Who Has Ever Been Afraid of Statistics and Analytics* can teach you the basic concepts that provide you with the knowledge to apply statistics and analytics in your life. You will also learn the most commonly used statistical methods and have the opportunity to practice those methods while using the Microsoft Excel spreadsheet program.

Please read the rest of this introduction so that you can become familiar with the distinctive features of this book. You can also visit the website for this book (www.ftpress.com/evenyoucanlearnstatistics3e) where you can learn more about this book as well as download files that support your learning of statistics.

Mathematics Is Always Optional!

Never mastered higher mathematics—or generally fearful of math? Not to worry, because in *Even You Can Learn Statistics and Analytics*, you will find that every concept is explained in plain English, without the use of higher mathematics or mathematical symbols. Interested in the mathematical foundations behind statistics? *Even You Can Learn Statistics and Analytics* includes **Equation Blackboards**, stand-alone sections that present the equations behind statistical methods and complement the main material. Either way, you can learn statistics.

Learning with the Concept-Interpretation Approach

Even You Can Learn Statistics and Analytics uses a **Concept-Interpretation** approach to help you learn statistics. For each important statistical concept, you will find the following:

- A **CONCEPT**, a plain language definition that uses no complicated mathematical terms.
- An **INTERPRETATION**, that fully explains the concept and its importance to statistics. When necessary, these sections also discuss common

misconceptions about the concept as well as the common errors people can make when trying to apply the concept.

For simpler concepts, an **EXAMPLES** section lists real-life examples or applications of the statistical concepts. For more involved concepts, **WORKED-OUT PROBLEMS** provide a complete solution to a statistical problem—including actual spreadsheet and calculator results—that illustrate how you can apply the concept to your own situations.

Practicing Statistics While You Learn Statistics

To help you learn statistics, you should always review the worked-out problems that appear in this book. As you review them, you can practice what you have just learned by using the optional **SPREADSHEET SOLUTION** sections.

Spreadsheet Solution sections enable you to use Microsoft Excel as you learn statistics.

Prefer to practice using a calculator? **CALCULATOR KEYS** sections (available online at www.ftpress.com/evenyoucanlearnstatistics3e) provide you with the step-by-step instructions to perform statistical analysis using one of the calculators from the Texas Instruments TI-83/84 family. (You can adapt many instruction sets for use with other TI statistical calculators.)

If you don't want to practice your spreadsheet skills, you can examine the spreadsheet results that appear throughout the book (or the calculator results available online). Many spreadsheet results are available as files that you can download for free at www.ftpress.com/evenyoucanlearnstatistics3e.

Spreadsheet program users will also benefit from Appendix D and Appendix E, which help teach you more about spreadsheets as you learn statistics.

And if technical issues or instructions have ever confounded your using Microsoft Excel in the past, check out Appendix A, which details the technical configuration issues you might face and explains the conventions used in all technical instructions that appear in this book.

important
point



In-Chapter Aids


As you read a chapter, look for the following icons for extra help:

Important Point icons highlight key definitions and explanations.



File icons identify files that allow you to examine the data in selected problems. (You can download these files for free at www.ftpress.com/evenyoucanlearnstatistics3e.)

interested
in
math?



Interested in the mathematical foundations of statistics? Then look for the Interested in Math? icons throughout the book. But remember, you can skip any or all of the math sections without losing any comprehension of the statistical methods presented, because math is always optional in this book!

End-of-Chapter Features

At the end of most chapters of *Even You Can Learn Statistics and Analytics*, you can find the following features, which you can review to reinforce your learning.

Important Equations

The **Important Equations** sections present all of the important equations discussed in the chapter. You can use these lists for reference and later study even if you have skipped over the Equation Blackboards and “interested in math” passages.

One-Minute Summaries

One-Minute Summaries are a quick review of the significant topics of a chapter in outline form. When appropriate, the summaries also help guide you to make the right decisions about applying statistics to the data you seek to analyze.

Test Yourself

The **Test Yourself** sections offer a set of short-answer questions and problems that enable you to review and test yourself (with answers provided) to see how much you have retained of the concepts presented in a chapter.

New to the Third Edition

The third edition of this book includes these features, which earlier editions did not contain:

- A new chapter that introduces the concepts and application of analytics, a new and growing part of statistics (Chapter 12)
- New chapters that illustrate descriptive and predictive analytical methods (Chapters 13 and 14)
- A new and expanded discussion about using Microsoft Excel, focused on using Excel 2013 (Microsoft Windows), Excel 2011 (OS X), and Office 365 Excel

The book also contains many revised in-chapter, worked-out problems and many new or revised end-of-chapter problems.

Summary

Even You Can Learn Statistics and Analytics can help you whether you are studying statistics as part of a formal course, just brushing up on your knowledge of statistics for a specific analysis, or need to learn about analytics. Be sure to visit the website for this book (www.ftpress.com/evenyoucanlearnstatistics3e) and feel free to contact the authors via email at davidlevine@davidlevinestatistics.com; include *Even You Can Learn Statistics and Analytics 3/e* in the subject line if you have any questions about this book.



Fundamentals of Statistics

- 1.1 The First Three Words of Statistics
 - 1.2 The Fourth and Fifth Words
 - 1.3 The Branches of Statistics
 - 1.4 Sources of Data
 - 1.5 Sampling Concepts
 - 1.6 Sample Selection Methods
- One-Minute Summary
Test Yourself

Every day, the media uses numbers to describe or analyze our world:

- “6 New Facts About Facebook” (A. Smith, www.pewresearch.org/author/asmith, 3 February 2014). A survey reported that women were more likely than men to cite seeing photos or videos, sharing with many people at once, seeing entertaining or funny posts, learning about ways to help others, and receiving support from people in their network as reasons to use Facebook.
- “First Two Years of College Wasted?” (M. Marklein, *USA Today*, 18 January 2011, p. 3A). A survey of more than 3,000 full-time, traditional-age students found that the students spent 51% of their time on socializing, recreation, and other activities; 9% of their time attending class and labs; and 7% of their time studying.
- “Follow the Tweets” (H. Rui, A. Whinston, and E. Winkler, *The Wall Street Journal*, 30 November 2009, p. R4). In this study, the authors found that the number of times a specific product was mentioned in comments in the Twitter social messaging service could be used to make accurate predictions of sales trends for that product.

You can make better sense of the numbers you encounter if you learn to understand statistics. **Statistics**, a branch of mathematics, uses procedures

that allow you to correctly analyze the numbers. These procedures, or **statistical methods**, transform numbers into useful information that you can use when making decisions about the numbers. Statistical methods can also tell you the known risks associated with making a decision as well as help you make more consistent judgments about the numbers.

Learning statistics requires you to reflect on the significance and the importance of the results to the decision-making process you face. This statistical interpretation means knowing when to ignore results because they are misleading, are produced by incorrect methods, or just restate the obvious, as in “100% of the authors of this book are named ‘David.’”

In this chapter, you begin by learning five basic words—*population*, *sample*, *variable*, *parameter*, and *statistic* (singular)—that identify the fundamental concepts of statistics. These five words, and the other concepts introduced in this chapter, help you explore and explain the statistical methods discussed in later chapters.

1.1 The First Three Words of Statistics

You’ve already learned that statistics is about analyzing things. Although *numbers* was the word used to represent things in the opening of this chapter, the first three words of statistics, *population*, *sample*, and *variable*, help you to better identify what you analyze with statistics.

Population

CONCEPT All the members of a group about which you want to reach a conclusion.

EXAMPLES All U.S. citizens who are currently registered to vote, all patients treated at a particular hospital last year, the entire set of individuals who accessed a website on a particular day.

Sample

CONCEPT The part of the population selected for analysis.

EXAMPLES The registered voters selected to participate in a recent survey concerning their intention to vote in the next election, the patients selected to fill out a patient satisfaction questionnaire, 100 boxes of cereal selected from a factory’s production line, 500 individuals who accessed a website on a particular day.

Variable

CONCEPT A characteristic of an item or an individual that will be analyzed using statistics.

EXAMPLES Gender, the party affiliation of a registered voter, the household income of the citizens who live in a specific geographical area, the publishing category (hardcover, trade paperback, mass-market paperback, textbook) of a book, the number of cell phones in a household.

INTERPRETATION All the variables taken together form the data of an analysis. Although people often say that they are analyzing their data, they are, more precisely, analyzing their variables.

You should distinguish between a variable, such as gender, and its value for an individual, such as male. An **observation** is all the values for an individual item in the sample. For example, a survey might contain two variables, gender and age. The first observation might be male, 40. The second observation might be female, 45. The third observation might be female, 55. By convention, when you organize data in tabular form, you place the values for a variable to be analyzed in a column. Therefore, some people refer to a variable as a *column of data*. Likewise, some people call an observation a *row of data*.

	Categorical Variables	Numerical Variables
Concept	The values of these variables are selected from an established list of categories.	The values of these variables involve a counted or measured value.
Subtypes	None	Discrete values are counts of things. Continuous values are measures and any value can theoretically occur, limited only by the precision of the measuring process.
Examples	Gender, a variable that has the categories “male” and “female.” variable. Academic major, a variable that might have the categories “English,” “Math,” “Science,” and “History,” among others.	The number of people living in a household, a discrete numerical variable. The time it takes for someone to commute to work, a continuous variable.



All variables should have an operational definition—that is, a universally accepted meaning that is understood by all associated with an analysis. Without operational definitions, confusion can occur. A famous example of such confusion was a survey that asked about sex and a number of survey takers answered yes and not male or female, as the survey writer intended.

1.2 The Fourth and Fifth Words

After you know what you are analyzing, or, using the words of Section 1.1, after you have identified the variables from the population or sample under study, you can define the **parameters** and **statistics** that your analysis will determine.

Parameter

CONCEPT A numerical measure that describes a variable (characteristic) from a population.

EXAMPLES The percentage of all registered voters who intend to vote in the next election, the percentage of all patients who are very satisfied with the care they received, the mean time that all visitors spent on a website during a particular day.

Statistic

CONCEPT A numerical measure that describes a variable (characteristic) of a sample (part of a population).

EXAMPLES The percentage of registered voters in a sample who intend to vote in the next election, the percentage of patients in a sample who are very satisfied with the care they received, the mean time that a sample of visitors spent on a website during a particular day.

INTERPRETATION Calculating statistics for a sample is the most common activity because collecting population data is impractical in many actual decision-making situations.

1.3 The Branches of Statistics

You can use parameters and statistics either to describe your variables or to reach conclusions about your data. These two uses define the two branches of statistics: **descriptive statistics** and **inferential statistics**.

Descriptive Statistics

CONCEPT The branch of statistics that focuses on collecting, summarizing, and presenting a set of data.

EXAMPLES The mean age of citizens who live in a certain geographical area, the mean length of all books about statistics, the variation in the time that visitors spent visiting a website.

INTERPRETATION You are most likely to be familiar with this branch of statistics because many examples arise in everyday life. Descriptive statistics serves as the basis for analysis and discussion in fields as diverse as securities trading, the social sciences, government, the health sciences, and professional sports. Descriptive methods can seem deceptively easy to apply because they are often easily accessible in calculating and computing devices. However, this ease does not mean that descriptive methods are without their pitfalls, as Chapter 2 and Chapter 3 explain.

Inferential Statistics

CONCEPT The branch of statistics that analyzes sample data to reach conclusions about a population.

EXAMPLE A survey that sampled 1,264 women found that 45% of those polled considered friends or family as their most trusted shopping advisers and only 7% considered advertising as their most trusted shopping adviser. By using methods discussed in Section 6.4, you can use these statistics to draw conclusions about the population of all women.

INTERPRETATION When you use inferential statistics, you start with a hypothesis and look to see whether the data are consistent with that hypothesis. This deeper level of analysis means that inferential statistical methods can be easily misapplied or misconstrued, and that many inferential methods require a calculating or computing device. (Chapters 6 through 9 discuss some of the inferential methods that you will most commonly encounter.)

1.4 Sources of Data

You begin every statistical analysis by identifying the source of the data that you will use for **data collection**. Among the important sources of data are **published sources**, **experiments**, and **surveys**.

Published Sources

CONCEPT Data available in print or in electronic form, including data found on Internet websites. Primary data sources are those published by the

individual or group that collected the data. Secondary data sources are those compiled from primary sources.

EXAMPLE Many U.S. federal agencies, including the Census Bureau, publish primary data sources that are available at the www.fedstats.gov website. Industry-specific groups and business news organizations commonly publish online or in-print secondary source data compiled by business organizations and government agencies.

INTERPRETATION You should always consider the possible bias of the publisher and whether the data contain all the necessary and relevant variables when using published sources. This is especially true of sources found through Internet search engines.

Experiments

CONCEPT A study that examines the effect on a variable of varying the value(s) of another variable or variables, while keeping all other things equal. A typical experiment contains both a treatment group and a control group. The treatment group consists of those individuals or things that receive the treatment(s) being studied. The control group consists of those individuals or things that do not receive the treatment(s) being studied.

EXAMPLE Pharmaceutical companies use experiments to determine whether a new drug is effective. A group of patients who have many similar characteristics is divided into two subgroups. Members of one group, the treatment group, receive the new drug. Members of the other group, the control group, often receive a placebo, a substance that has no medical effect. After a time period, statistics about each group are compared.

INTERPRETATION Proper experiments are either single-blind or double-blind. A study is a single-blind experiment if only the researcher conducting the study knows the identities of the members of the treatment and control groups. If neither the researcher nor study participants know who is in the treatment group and who is in the control group, the study is a double-blind experiment.

When conducting experiments that involve placebos, researchers also have to consider the placebo effect—that is, whether people in the control group will improve because they believe they are getting a real substance that is intended to produce a positive result. When a control group shows as much improvement as the treatment group, a researcher can conclude that the placebo effect is a significant factor in the improvements of both groups.

Surveys

CONCEPT A process that uses questionnaires or similar means to gather values for the responses from a set of participants.

EXAMPLES The decennial U.S. census mail-in form, a poll of likely voters, a website instant poll or “question of the day.”

INTERPRETATION Surveys are either **informal**, open to anyone who wants to participate; **targeted**, directed toward a specific group of individuals; or include people chosen at random. The type of survey affects how the data collected can be used and interpreted.

1.5 Sampling Concepts

In the definition of **statistic** in Section 1.2, you learned that calculating statistics for a sample is the most common activity because collecting population data is usually impractical. Because samples are so commonly used, you need to learn the concepts that help identify all the members of a population and that describe how samples are formed.

Frame

CONCEPT The list of all items in the population from which the sample will be selected.

EXAMPLES Voter registration lists, municipal real estate records, customer or human resource databases, directories.

INTERPRETATION Frames influence the results of an analysis, and using different frames can lead to different conclusions. You should always be careful to make sure your frame completely represents a population; otherwise, any sample selected will be biased, and the results generated by analyses of that sample will be inaccurate.

Sampling

CONCEPT The process by which members of a *population* are selected for a *sample*.

EXAMPLES Choosing every fifth voter who leaves a polling place to interview, selecting playing cards randomly from a deck, polling every tenth visitor who views a certain website today.

INTERPRETATION Some sampling techniques, such as an “instant poll” found on a web page, are naturally suspect as such techniques do not depend on a well-defined frame. The sampling technique that uses a well-defined frame is **probability sampling**.

Probability Sampling

CONCEPT A sampling process that considers the chance of selection of each item. Probability sampling increases your chance that the sample will be representative of the population.

EXAMPLES The registered voters selected to participate in a recent survey concerning their intention to vote in the next election, the patients selected to fill out a patient-satisfaction questionnaire, 100 boxes of cereal selected from a factory’s production line.

INTERPRETATION You should use probability sampling whenever possible, because *only* this type of sampling enables you to apply inferential statistical methods to the data you collect. In contrast, you should use nonprobability sampling, in which the chance of occurrence of each item being selected is not known, to obtain rough approximations of results at low cost or for small-scale, initial, or pilot studies that will later be followed up by a more rigorous analysis. Surveys and polls that invite the public to call in or answer questions on a web page are examples of nonprobability sampling.

Simple Random Sampling

CONCEPT The probability sampling process in which every individual or item from a population has the same chance of selection as every other individual or item. Every possible sample of a certain size has the same chance of being selected as every other sample of that size.

EXAMPLES Selecting a playing card from a shuffled deck or using a statistical device such as a table of random numbers.

INTERPRETATION Simple random sampling forms the basis for other random sampling techniques. The word *random* in this phrase requires clarification. In this phrase, *random* means no repeating patterns—that is, in a given sequence, a given pattern is equally likely (or unlikely). It does not refer to the most commonly used meaning of “unexpected” or “unanticipated” (as in “random acts of kindness”).

Other Probability Sampling Methods

Other, more complex, sampling methods are also used in survey sampling. In a stratified sample, the items in the frame are first subdivided into separate subpopulations, or strata, and a simple random sample is selected within each of the strata. In a cluster sample, the items in the frame are divided into several clusters so that each cluster is representative of the entire population. A random sampling of clusters is then taken, and all the items in each selected cluster or a sample from each cluster are then studied.

1.6 Sample Selection Methods

Sampling can be done either with or without replacement of the items being selected. Almost all survey sampling is done without replacement.

Sampling with Replacement

CONCEPT A sampling method in which each selected item is returned to the frame from which it was selected so that it has the same probability of being selected again.

EXAMPLE Selecting items from a fishbowl and returning each item to it after the selection is made.

Sampling Without Replacement

CONCEPT A sampling method in which each selected item is not returned to the frame from which it was selected. Using this technique, an item can be selected no more than one time.

EXAMPLES Selecting numbers in state lottery games, selecting cards from a deck of cards during games of chance such as blackjack or poker.

INTERPRETATION Sampling without replacement means that an item can be selected no more than one time. You should choose sampling without replacement instead of sampling with replacement because statisticians generally consider the former to produce more desirable samples.



spreadsheet solution

Entering Data

Enter the data values of a variable in a blank column of a worksheet. Use the row 1 cell for the variable name.

To create a new file (workbook) that contains a blank worksheet for your entries, select **File** → **New** and, in the **New** panel, double-click the **Blank Workbook** icon. (In Excel 2007, click the **Office Button** instead of selecting **File**.) To enter data into a specific cell, move the cell pointer to that cell by using the cursor keys, moving the mouse pointer, or completing the proper touch operation. As you type an entry, the entry appears in the formula bar area located over the top of the worksheet. You complete your entry by pressing **Tab** or **Enter** or by clicking the checkmark button in the formula bar.

To save your new file, select **File** → **Save As** and, in the **Save As** dialog box, navigate to the folder where you want to save your file. Accept or revise the filename and then click **Save**. To later retrieve the file, select **File** → **Open** and in the **Open** dialog box, navigate to the folder that contains the desired file, select the desired file from the list, and then click **Open**. (In Excel 2007, you begin these operations by clicking the **Office Button**, not by selecting **File**.)



Throughout this book, the symbol → links a sequence of Ribbon or menu selections. **File** → **New** means to first select the **File** tab and then select **New** from the list that appears.

One-Minute Summary

Mastering basic vocabulary is the first step in learning statistics. Understanding the types of statistical methods, the sources of data used for data collection, sampling methods, and the types of variables used in statistical analysis are also important introductory concepts. Subsequent chapters focus on four important reasons for learning statistics:

- To present and describe information (Chapters 2 and 3)
- To reach conclusions about populations based only on sample results (Chapters 4 through 9)
- To develop reliable forecasts (Chapters 10 and 11)
- To use analytics to reach conclusions about large sets of data (Chapters 12, 13, 14)

Test Yourself

1. The portion of the population that is selected for analysis is called:
 - (a) a sample
 - (b) a frame
 - (c) a parameter
 - (d) a statistic

2. A summary measure that is computed from only a sample of the population is called:
 - (a) a parameter
 - (b) a population
 - (c) a discrete variable
 - (d) a statistic

3. The height of an individual is an example of a:
 - (a) discrete variable
 - (b) continuous variable
 - (c) categorical variable
 - (d) constant

4. The body style of an automobile (sedan, minivan, SUV, and so on) is an example of a:
 - (a) discrete variable
 - (b) continuous variable
 - (c) categorical variable
 - (d) constant

5. The number of credit cards in a person's wallet is an example of a:
 - (a) discrete variable
 - (b) continuous variable
 - (c) categorical variable
 - (d) constant

6. Statistical inference occurs when you:
 - (a) compute descriptive statistics from a sample
 - (b) take a complete census of a population
 - (c) present a graph of data
 - (d) take the results of a sample and reach conclusions about a population

7. The human resources director of a large corporation wants to develop a dental benefits package and decides to select 100 employees from a list of all 5,000 workers in order to study their preferences for the various components of a potential package. All the employees in the corporation constitute the _____.
 - (a) sample
 - (b) population
 - (c) statistic
 - (d) parameter
8. The human resources director of a large corporation wants to develop a dental benefits package and decides to select 100 employees from a list of all 5,000 workers in order to study their preferences for the various components of a potential package. The 100 employees who will participate in this study constitute the _____.
 - (a) sample
 - (b) population
 - (c) statistic
 - (d) parameter
9. Those methods that involve collecting, presenting, and computing characteristics of a set of data in order to properly describe the various features of the data are called:
 - (a) statistical inference
 - (b) the scientific method
 - (c) sampling
 - (d) descriptive statistics
10. Based on the results of a poll of 500 registered voters, the conclusion that the Democratic candidate for U.S. president will win the upcoming election is an example of:
 - (a) inferential statistics
 - (b) descriptive statistics
 - (c) a parameter
 - (d) a statistic
11. A numerical measure that is computed to describe a characteristic of an entire population is called:
 - (a) a parameter
 - (b) a population
 - (c) a discrete variable
 - (d) a statistic

12. You were working on a project to examine the value of the American dollar as compared to the English pound. You accessed an Internet site where you obtained this information for the past 50 years. Which method of data collection were you using?
 - (a) published sources
 - (b) experimentation
 - (c) surveying
13. Which of the following is a discrete variable?
 - (a) The favorite flavor of ice cream of students at your local elementary school
 - (b) The time it takes for a certain student to walk to your local elementary school
 - (c) The distance between the home of a certain student and the local elementary school
 - (d) The number of teachers employed at your local elementary school
14. Which of the following is a continuous variable?
 - (a) The eye color of children eating at a fast-food chain
 - (b) The number of employees of a branch of a fast-food chain
 - (c) The temperature at which a hamburger is cooked at a branch of a fast-food chain
 - (d) The number of hamburgers sold in a day at a branch of a fast-food chain
15. The number of cell phones in a household is an example of:
 - (a) a categorical variable
 - (b) a discrete variable
 - (c) a continuous variable
 - (d) a statistic

Answer True or False:

16. The possible responses to the question, "How long have you been living at your current residence?" are values from a continuous variable.
17. The possible responses to the question, "How many times in the past seven days have you streamed a movie or TV show online?" are values from a discrete variable.

Fill in the blank:

18. An insurance company evaluates many variables about a person before deciding on an appropriate rate for automobile insurance. The number of accidents a person has had in the past three years is an example of a _____ variable.

19. An insurance company evaluates many variables about a person before deciding on an appropriate rate for automobile insurance. The distance a person drives in a day is an example of a _____ variable.
20. An insurance company evaluates many variables about a person before deciding on an appropriate rate for automobile insurance. A person's marital status is an example of a _____ variable.
21. A numerical measure that is computed from only a sample of the population is called a _____.
22. The portion of the population that is selected for analysis is called the _____.
23. A college admission application includes many variables. The number of advanced placement courses the student has taken is an example of a _____ variable.
24. A college admission application includes many variables. The gender of the student is an example of a _____ variable.
25. A college admission application includes many variables. The distance from the student's home to the college is an example of a _____ variable.

Answers to Test Yourself

- | | |
|-------|-----------------|
| 1. a | 14. c |
| 2. d | 15. b |
| 3. b | 16. True |
| 4. c | 17. True |
| 5. a | 18. discrete |
| 6. d | 19. continuous |
| 7. b | 20. categorical |
| 8. a | 21. statistic |
| 9. d | 22. sample |
| 10. a | 23. discrete |
| 11. a | 24. categorical |
| 12. a | 25. continuous |
| 13. d | |

References

1. Berenson, M. L., D. M. Levine, and K. A. Szabat. *Basic Business Statistics: Concepts and Applications*, Thirteenth Edition. Upper Saddle River, NJ: Pearson Education, 2015.
2. Cochran, W. G. *Sampling Techniques*, Third Edition. New York: John Wiley & Sons, 1977.



Index

A

α , 141

Alternative hypothesis,
137–139, 349

Analysis of variance
(ANOVA). *See* One-Way
Analysis of Variance

ANOVA summary table,
184–185

Analytics, 257–258, 349

Arithmetic mean, 37–38,
352

Arithmetic and algebra
review, 295–304

B

β , 141

Bad charts, 28–29

Bar chart, 16–17, 349

Big Data, 258, 349

Binomial distribution,
86–89, 349

Bootstrapping, 126–128

Box-and-whisker plot,
52–56, 349

Bullet graphs, 271–272,
350

C

Categorical variable, 3, 350

Cell, 350

Central limit theorem,
114–116

Certain event, 69

Chi-square distribution,
178, 350

Chi-square distribution
tables, 314–315

Chi-square test, 175–182

Classical approach to
probability, 73

Classification trees,
279–281, 350

Cluster analysis, 283–285,
350

Coefficient of correlation,
218, 350

Coefficient of
determination, 217–218,
350

Coefficient of multiple
determination, 242–243

Collectively exhaustive
events, 69–70, 350

Complement, 70

Completely randomized
design. *See* One-Way
Analysis of Variance

Confidence interval
estimate, 118–120, 350
For the mean
(σ unknown),
120–123
For the proportion,
123–126
For the slope, 221–222,
245–246

Continuous numerical
variables, 3, 350

Critical value, 141, 350

D

Dashboards, 265–266, 350

Data mining, 260, 350

Degrees of freedom, 161,
178, 350

Dependent variable, 204,
350

Descriptive analytics, 259,
351

Descriptive statistics, 4–5,
351

- Discrete numerical variables, 3, 351
- Discrete probability distribution, 80–81
- Discrete values, 3
- Double-blind study, 6
- Downloadable files, 345–348
- Drilldown, 266–267, 351

- E**
- Elementary event, 69
- Empirical approach to probability, 73–74
- Equation Blackboard
 - Binomial distribution, 87–88
 - Chi-square tests, 182
 - Confidence interval estimate for the mean (σ unknown), 122
 - Confidence interval estimate for the proportion, 125
 - Confidence interval estimate for the slope, 222
 - Mean, 39
 - Mean and standard deviation of a discrete probability distribution, 85
 - Median, 41
 - Paired t test, 165–166
 - One Way Analysis of Variance (ANOVA), 187–189
 - Poisson distribution, 91
 - Pooled-variance t test for the difference in two means, 160–161
 - Range, 46
 - Regression measures of variation, 216–217
 - Slope, 209–212
 - Standard deviation, 49
 - Standard error of the estimate, 219
 - t test for the slope, 220–221
 - Variance, 49
 - Y intercept, 209–212
 - Z scores, 50
 - Z test for the difference in two proportions, 154–155
- Event, 68, 351
- Expected frequency, 176, 351

Expected value of a
variable, 81–82, 351
Experiments, 6, 351
Explanatory variable, 204,
351

F–G

F distribution 184, 351
F distribution tables,
316–331
F test statistic, 184
Factor, 180
Five number summary, 52,
351
Frame, 7, 351
Frequency distribution,
22–24, 338–339, 351

H

Histogram, 24–25, 351
Hypothesis testing, 351
Hypothesis testing steps,
143

I–J–K

Independent events, 72,
351

Independent variable, 204,
352
Inferential statistics, 5, 352
Joint event, 69, 352

L

Least-squares method,
206–207
Left-skewed, 51
Level of significance, 141,
352

M

Mean, 37–38, 352
Mean Squares, 352
 Among Groups (*MSA*),
 184–185
 Total (*MST*), 184–185
 Within Groups (*MSW*),
 184–185
Measures of
 Central tendency, 37–41
 Position, 41–45
 Variation, 45–50
Median, 38–41, 352
Mode, 41, 352
Multidimensional scaling,
286–288, 352

Multiple regression model,
239, 352
Mutually exclusive, 70, 352

N

Net regression coefficients,
241
Normal distribution,
93–99, 352
Normal distribution tables,
306–309
Normal probability plot,
100–101, 352
Null event, 69
Null hypothesis, 138, 352
Numerical variable, 3, 352

O

Observation, 3
Observed frequency, 352
One-tail test, 144–145
One-Way Analysis of
Variance, 182–192
Assumptions, 192
Operational definition, 3
Overall *F* test, 243

P

p-value, 143–144, 352
Paired *t* test, 162–167, 352
Parameter, 4, 352
Pareto chart, 18–19, 353
Percentage distribution,
22–24, 353
Pie chart, 17–18, 353
PivotTables, 337–338
Placebo, 6
Point estimate, 117
Poisson distribution,
89–92, 353
Pooled-variance *t* test,
156–162
Population, 2, 353
Power of the test, 142
Practical significance, 141
Prediction, 203–204
Predictive analytics,
259–260
Primary data sources, 5
Probability, 69, 353
Rules, 70–73
Probability distribution for
a continuous variable,
92

Probability distribution for
discrete random
variables, 79, 353
Probability sampling, 8, 353
Published sources, 5, 353

Q–R

Quartiles, 41–44
Random variable, 69–70
Range, 46, 353
Region of nonrejection,
140
Region of rejection, 140,
353
Regression model
prediction, 209
Regression analysis,
203–204, 353
Regression assumptions,
212–215
Regression trees, 281–283,
353
Residual analysis, 354
Simple linear
regression, 213–214
Multiple regression, 244
Response variable, 204,
350, 354
Right-skewed, 51

S

Sample, 2, 354
Sampling, 7, 354
Sampling distribution, 114,
354
Of the mean, 115–116
Of the proportion,
116–117
Sampling error, 117–118,
354
Sampling with
replacement, 9, 354
Sampling without
replacement, 9, 354
Scatter plot, 26–27,
204–205, 354
Shape, 51–56
Simple linear regression,
205–213, 354
Assumptions, 212–213
Simple random sampling,
8, 354
Skewness, 354
Slope, 206, 354
Software for analytics,
260–261
Sparklines, 268–269, 354
Spreadsheet operating
conventions, 293

- Spreadsheet Solutions
 - Bar and pie charts, 18
 - Binomial probabilities, 88
 - Bullet graphs, 272
 - Chi-square tests, 181
 - Classification and regression trees, 283
 - Confidence interval estimate for the mean (σ unknown), 123, 339–340
 - Confidence interval estimate for the proportion, 126
 - Cluster analysis, 285
 - Drill-down, 268
 - Entering data, 10
 - Frequency distributions and histograms, 25
 - Measures of central tendency and position, 45
 - Measures of shape, 56
 - Measures of variation, 48
 - Multidimensional scaling, 288
 - Multiple regression, 246
 - Normal probabilities, 99
 - One-Way Analysis of Variance (ANOVA), 189–190
 - Paired t test, 164, 340
 - Pareto charts, 20
 - Poisson probabilities, 92
 - Pooled-variance t test for the difference in two means, 158, 160
 - Scatter plots, 28
 - Simple linear regression, 222, 340–342
 - Sparklines, 269
 - Time-series plots, 28
 - Treemap, 271
 - Two-way tables, 22
 - Z test for the difference in two proportions, 154
- Spreadsheet technical configurations, 294
- Spreadsheet Tips
 - Analysis ToolPak Tips, 342–344

- Chart Tips, 333–334
- Function tips, 335
- Standard deviation, 46–48, 354
- Standard deviation of a variable, 82–85
- Standard error of the estimate, 219, 354
- Standard (Z) scores, 50
- Statistic, 4, 354
- Statistics, 1, 355
- Subjective approach to probability, 74
- Sum of Squares
 - Error (SSE), 215, 351
 - Regression (SSR), 215, 353
 - Total (SST), 183, 215, 355
 - Among Groups (SSA), 183, 355
 - Within Groups (SSW), 184, 355
- Summary table, 15–16, 355
- Surveys, 7, 355
- Symmetric, 51, 355

T–U

- t distribution, 120, 355
- t distribution tables, 310–313
- Tables of the
 - Chi-square distribution, 314–315
 - F distribution, 316–331
 - Normal distribution, 306–309
 - t distribution, 310–313
- Test of hypothesis
 - Chi-square test, 175–182
 - For the difference between two proportions, 149–155
 - For the difference between the means of two independent groups, 156–162
 - For the slope, 220
 - In multiple regression, 245
 - One-Way Analysis of Variance, 182–193
 - Paired t test, 162–167

Test statistic, 139–140, 355

Time-series plot, 25–26,
355

Treatment effect, 183

Treemap, 270, 355

Two-tail test, 144–145

Two-way cross-
classification tables,
20–22, 355

Type I error, 141, 355

Type II error, 141–142, 356

V–W

Variable, 3, 356

Variance, 46–48, 356

X–Y–Z

Y intercept, 206, 356

Z scores, 50, 356