

Excel for Statistics

Thomas J. Quirk

Excel 2019 for Business Statistics

A Guide to Solving Practical Problems

Second Edition

 Springer

Excel for Statistics

Excel for Statistics is a series of textbooks that explain how to use Excel to solve statistics problems in various fields of study. Professors, students, and practitioners will find these books teach how to make Excel work best in their respective fields. Applications include any disciplines that use data and can benefit from the power and simplicity of Excel. Books cover all the steps for running statistical analyses in Excel 2019, Excel 2016, and Excel 2013. The approach also teaches critical statistics skills, making the books particularly applicable for statistics courses taught outside of mathematics or statistics departments.

Series editor: Thomas J. Quirk

The following books are in this series:

T.J. Quirk, E. Rhiney, *Excel 2019 for Advertising Statistics: A Guide to Solving Practical Problems*, Excel for Statistics, Second Edition. Springer International Publishing Switzerland 2020.

T.J. Quirk, E. Rhiney, *Excel 2019 for Marketing Statistics: A Guide to Solving Practical Problems*, Excel for Statistics, Second Edition. Springer International Publishing Switzerland 2020.

T.J. Quirk, M. Quirk, H.F. Horton, *Excel 2019 for Biological and Life Sciences Statistics: A Guide to Solving Practical Problems*, Excel for Statistics, Second Edition. Springer International Publishing AG, part of Springer Nature 2020.

T.J. Quirk, *Excel 2019 for Business Statistics: A Guide to Solving Practical Problems*, Excel for Statistics, Second Edition. Springer International Publishing AG, part of Springer Nature 2020.

T.J. Quirk, *Excel 2019 for Engineering Statistics: A Guide to Solving Practical Problems*, Excel for Statistics, Second Edition. Springer International Publishing AG, part of Springer Nature 2020.

T.J. Quirk, *Excel 2019 for Educational and Psychological Statistics: A Guide to Solving Practical Problems*, Excel for Statistics, Second Edition. Springer International Publishing AG, part of Springer Nature 2020.

T.J. Quirk, *Excel 2019 for Social Science Statistics: A Guide to Solving Practical Problems*, Excel for Statistics, Second Edition. Springer International Publishing AG, part of Springer Nature 2020.

T.J. Quirk, *Excel 2016 Applied Statistics for High School Students: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2018.

T.J. Quirk, E. Rhiney, *Excel 2016 for Advertising Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2017.

T.J. Quirk, S. Cummings, *Excel 2016 for Social Work Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2017.

T.J. Quirk, E. Rhiney, *Excel 2016 for Marketing Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, *Excel 2016 for Business Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk. *Excel 2016 for Engineering Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, M. Quirk, H.F. Horton, *Excel 2016 for Biological and Life Sciences Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk. *Excel 2016 for Educational and Psychological Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, *Excel 2016 for Social Science Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, M. Quirk, H. Horton, *Excel 2016 for Physical Sciences Statistics: A Guide to Solving Practical Problems*. Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, S. Cummings, *Excel 2016 for Health Services Management Statistics: A Guide to Solving Practical Problems*. Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, J. Palmer-Schuyler, *Excel 2016 for Human Resource Management Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, M. Quirk, H.F. Horton. *Excel 2016 for Environmental Sciences Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, M. Quirk, H.F. Horton. *Excel 2013 for Physical Sciences Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, S. Cummings, *Excel 2013 for Health Services Management Statistics: A Guide to Solving Practical Problems*. Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, J. Palmer-Schuyler, *Excel 2013 for Human Resource Management Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2016.

T.J. Quirk, *Excel 2013 for Business Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2015.

T.J. Quirk. *Excel 2013 for Engineering Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2015.

T.J. Quirk, M. Quirk, H.F. Horton, *Excel 2013 for Biological and Life Sciences Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2015.

T.J. Quirk. *Excel 2013 for Educational and Psychological Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2015.

T.J. Quirk, *Excel 2013 for Social Science Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2015.

T.J. Quirk, M. Quirk, H.F. Horton, *Excel 2013 for Environmental Sciences Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2015.

T.J. Quirk, M. Quirk, H.F. Horton, *Excel 2010 for Environmental Sciences Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2015.

T.J. Quirk, J. Palmer-Schuyler, *Excel 2010 for Human Resource Management Statistics: A Guide to Solving Practical Problems*, Excel for Statistics. Springer International Publishing Switzerland 2014.

Additional Statistics books by Dr. Thomas J. Quirk that have been published by Springer

T.J. Quirk, *Excel 2010 for Business Statistics: A Guide to Solving Practical Problems*. Springer Science+Business Media 2011.

T.J. Quirk, *Excel 2010 for Engineering Statistics: A Guide to Solving Practical Problems*. Springer International Publishing Switzerland 2014.

T.J. Quirk, S. Cummings, *Excel 2010 for Health Services Management Statistics: A Guide to Solving Practical Problems*. Springer International Publishing Switzerland 2014.

T.J. Quirk, M. Quirk, H. Horton, *Excel 2010 for Physical Sciences Statistics: A Guide to Solving Practical Problems*. Springer International Publishing Switzerland 2013.

T.J. Quirk, M. Quirk, H.F. Horton, *Excel 2010 for Biological and Life Sciences Statistics: A Guide to Solving Practical Problems*. Springer Science+Business Media New York 2013.

T.J. Quirk, *Excel 2010 for Social Science Statistics: A Guide to Solving Practical Problems*. Springer Science+Business Media New York 2012.

T.J. Quirk, *Excel 2010 for Educational and Psychological Statistics: A Guide to Solving Practical Problems*. Springer Science+Business Media New York 2012.

T.J. Quirk, *Excel 2007 for Business Statistics: A Guide to Solving Practical Problems*. Springer Science+Business Media New York 2012.

T.J. Quirk, *Excel 2007 for Educational and Psychological Statistics: A Guide to Solving Practical Problems*. Springer Science+Business Media New York 2012.

T.J. Quirk, *Excel 2007 for Social Science Statistics: A Guide to Solving Practical Problems*. Springer Science+Business Media New York 2012.

T.J. Quirk, *Excel 2007 for Biological and Life Sciences Statistics: A Guide to Solving Practical Problems*. Springer Science+Business Media New York 2013.

More information about this series at <http://www.springer.com/series/13491>

Thomas J. Quirk

Excel 2019 for Business Statistics

A Guide to Solving Practical Problems

Second Edition

 Springer

Thomas J. Quirk
Professor of Marketing
Webster University
St. Louis, MO, USA

ISSN 2570-4605

ISSN 2570-4613 (electronic)

Excel for Statistics

ISBN 978-3-030-39260-4

ISBN 978-3-030-39261-1 (eBook)

<https://doi.org/10.1007/978-3-030-39261-1>

© Springer Nature Switzerland AG 2016, 2020

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

This book is dedicated to more than 3000 students I have taught at Webster University's campuses in St. Louis, London, and Vienna; the students at Principia College in Elsau, Illinois; and the students at the Cooperative State University of Baden-Wuerttemberg in Heidenheim, Germany. These students taught me a great deal about the art of teaching. I salute them all, and I thank them for helping me to become a better teacher.

Thomas J. Quirk

Preface

Excel 2019 for Business Statistics: A Guide to Solving Practical Problems updates the Excel steps and screenshots from the previously published *Excel 2016 for Business Statistics: A Guide to Solving Practical Problems*, and it contains a number of important changes. The explanations of statistics and statistical formulas have been made clearer. The Excel steps now match perfectly the Excel 2019 version. Thirty percent of the end-of-chapter problems, and their answers in an Appendix, are new to this book. Thirty percent of the 160+ screenshots are new so that they match the new Excel commands to ensure that you are using Excel correctly each step of the way.

The eight chapters in the book (Mean, Standard Deviation, and Standard Error of the Mean; Random Sampling; Confidence Interval about the Mean; One-Group t-test; Two Group t-test; Correlation and Linear Regression; Multiple Correlation; and One-Way Analysis of Variance) have been rewritten to improve their explanation of statistics. The answers to all of the problems in the book are provided, and there is a Practice Test so that you can test your ability to solve statistics problems using Excel. This book is an introduction to statistics, not a full-blown explanation of statistics.

A word of caution: This book does not attempt to teach you all of the “bells and whistles” of Excel 2019. We have left that objective to other books. Instead, this book will teach you the Excel steps you need to solve the interesting problems in the book. You should think of Excel as merely the “computer language” needed to solve statistics problems. In a sense, this approach is similar to the one you would need if you planned to spend a year living in Europe in Vienna, Austria, where you needed to learn some basic German (e.g., “How much does this cost?” “Where is the train station?” “Please give me the bill for my dinner” “How can I get to the airport?”) but you do not need to become fluent in that language to survive. This book focuses on the Excel steps needed to solve the problems in the book. The task of showing you how to use the many powers of Excel are beyond the scope of this book.

This book was written by a Professor who wanted to respond to the complaints of many students about their inability to understand their statistics textbook and about

their inability to understand their professor's explanation of theoretical statistics. This book is self-instructional and does not depend on a professor's explanation of statistics. This book will teach you the general concepts of statistics without burying you in dull statistical theory. You will learn *why* you are performing the Excel steps through the objectives included in the chapters. The statistical concepts and practice problems get progressively more sophisticated as they build on what you have already learned from studying this book. This book is understandable by both undergraduate and graduate students who are taking their first course in statistics, by researchers, and by working professionals who want to solve interesting problems in their chosen field of study.

This book was written by a Professor who is, first and foremost, committed to helping you to understand how to use statistics to solve interesting problems in your chosen field of study. The ideas in this book have been classroom tested over the past 11 years in both undergraduate and graduate courses at Webster University, a liberal arts college located in St. Louis, Missouri, in the middle of the USA. This book is part of a series of more than 30 introductory statistics textbooks, in 12 fields of study, that have been published by Springer by Prof. Quirk, which have helped thousands of students, researchers, and working professionals learn how to use Excel to solve interesting statistics problems. These fields of study include: (1) Business, (2) Education/Psychology, (3) Social Science, (4) Biological and Life Sciences, (5) Physical Sciences, (6) Engineering, (7) Health Services Management, (8) Human Resource Management, (9) Environmental Sciences, (10) Marketing, (11) Social Work, and (12) Advertising.

Prof. Thomas J. Quirk is Professor Emeritus of Marketing at The Walker School of Business and Technology at Webster University in St. Louis, Missouri (US) where he taught Marketing Statistics, Marketing Research, and Pricing Strategies. At the beginning of his academic career, Prof. Quirk spent 6 years in educational research at The American Institutes for Research and Educational Testing Service. Prof. Quirk has published more than 20 articles in professional journals including *The Journal of Educational Psychology*, *Journal of Educational Research*, *Review of Educational Research*, *Journal of Educational Measurement*, and *Educational Technology*, published more than 60 textbook supplements in Management and Marketing, and presented more than 20 papers at professional meetings, including annual meetings of the American Educational Research Association, the American Psychological Association, and the National Council on Measurement in Education. Prof. Quirk holds a BS in Mathematics from John Carroll University, both an MA in Education and a PhD in Educational Psychology from Stanford University, and an MBA from the University of Missouri-St. Louis.

St. Louis, MO, USA

Thomas J. Quirk

Acknowledgements

Excel 2019 for Business Statistics: A Guide to Solving Practical Problems is the result of inspiration from three important people: my two daughters and my wife. Jennifer Quirk McLaughlin invited me to visit her MBA classes several times at the University of Witwatersrand in Johannesburg, South Africa. These visits to a first-rate MBA program convinced me that there was a need for a book to teach students how to solve practical problems using Excel. Meghan Quirk-Horton's dogged dedication to learning the many statistical techniques needed to complete her PhD dissertation illustrated the need for a statistics book that would make this daunting task more user-friendly. And Lynne Buckley-Quirk was the number-one cheerleader for this project from the beginning, always encouraging me and helping me remain dedicated to completing it.

Thomas J. Quirk

Contents

1	Sample Size, Mean, Standard Deviation, and Standard Error of the Mean	1
1.1	Mean	1
1.2	Standard Deviation	2
1.3	Standard Error of the Mean	3
1.4	Sample Size, Mean, Standard Deviation, and Standard Error of the Mean	4
1.4.1	Using the Fill/Series/Columns Commands	4
1.4.2	Changing the Width of a Column	6
1.4.3	Centering Information in a Range of Cells	6
1.4.4	Naming a Range of Cells	8
1.4.5	Finding the Sample Size Using the =COUNT Function	9
1.4.6	Finding the Mean Score Using the =AVERAGE Function	10
1.4.7	Finding the Standard Deviation Using the =STDEV Function	10
1.4.8	Finding the Standard Error of the Mean	10
1.5	Saving a Spreadsheet	13
1.6	Printing a Spreadsheet	14
1.7	Formatting Numbers in Currency Format (Two Decimal Places)	15
1.8	Formatting Numbers in Number Format (Three Decimal Places)	17
1.9	End-of-Chapter Practice Problems	17
	Reference	20
2	Random Number Generator	21
2.1	Creating Frame Numbers for Generating Random Numbers	21
2.2	Creating Random Numbers in an Excel Worksheet	25

- 2.3 Sorting Frame Numbers into a Random Sequence 26
- 2.4 Printing an Excel File So That All of the Information
Fits onto One Page 29
- 2.5 End-of-Chapter Practice Problems 34
- 3 Confidence Interval About the Mean Using the TINV
Function and Hypothesis Testing 37**
 - 3.1 Confidence Interval About the Mean 37
 - 3.1.1 How to Estimate the Population Mean 37
 - 3.1.2 Estimating the Lower Limit and the Upper Limit
of the 95% Confidence Interval About the Mean 38
 - 3.1.3 Estimating the Confidence Interval for the Chevy
Impala in Miles Per Gallon 39
 - 3.1.4 Where Did the Number “1.96” Come From? 40
 - 3.1.5 Finding the Value for t in the Confidence
Interval Formula 41
 - 3.1.6 Using Excel’s TINV Function to Find the Confidence
Interval About the Mean 42
 - 3.1.7 Using Excel to Find the 95% Confidence Interval
for a Car’s mpg Claim 42
 - 3.2 Hypothesis Testing 48
 - 3.2.1 Hypotheses Always Refer to the Population
of People or Events That You Are Studying 49
 - 3.2.2 The Null Hypothesis and the Research
(Alternative) Hypothesis 50
 - 3.2.3 The Seven Steps for Hypothesis-Testing
Using the Confidence Interval About the Mean 53
 - 3.3 Alternative Ways to Summarize the Result
of a Hypothesis Test 59
 - 3.3.1 Different Ways to Accept the Null Hypothesis 59
 - 3.3.2 Different Ways to Reject the Null Hypothesis 60
 - 3.4 End-of-Chapter Practice Problems 60
 - References 66
- 4 One-Group t-Test for the Mean 67**
 - 4.1 The Seven STEPS for Hypothesis-Testing Using
the One-Group t-Test 67
 - 4.1.1 STEP 1: State the Null Hypothesis
and the Research Hypothesis 68
 - 4.1.2 STEP 2: Select the Appropriate Statistical Test 68
 - 4.1.3 STEP 3: Decide on a Decision Rule
for the One-Group t-Test 68
 - 4.1.4 STEP 4: Calculate the Formula
for the One-Group t-Test 69
 - 4.1.5 STEP 5: Find the Critical Value of t in the t-Table
in Appendix E 70

- 4.1.6 STEP 6: State the Result of Your Statistical Test 71
- 4.1.7 STEP 7: State the Conclusion of Your Statistical Test in Plain English! 71
- 4.2 One-Group t-Test for the Mean 72
- 4.3 Can You Use Either the 95% Confidence Interval About the Mean OR the One-Group t-Test When Testing Hypotheses? 77
- 4.4 End-of-Chapter Practice Problems 77
- References 82
- 5 Two-Group t-Test of the Difference of the Means for Independent Groups 83**
 - 5.1 The Nine STEPS for Hypothesis-Testing Using the Two-Group t-Test 84
 - 5.1.1 STEP 1: Name One Group, Group 1, and the Other Group, Group 2 84
 - 5.1.2 STEP 2: Create a Table That Summarizes the Sample Size, Mean Score, and Standard Deviation of Each Group 84
 - 5.1.3 STEP 3: State the Null Hypothesis and the Research Hypothesis for the Two-Group t-Test 86
 - 5.1.4 STEP 4: Select the Appropriate Statistical Test 86
 - 5.1.5 STEP 5: Decide on a Decision Rule for the Two-Group t-Test 86
 - 5.1.6 STEP 6: Calculate the Formula for the Two-Group t-Test 86
 - 5.1.7 STEP 7: Find the Critical Value of t in the t-Table in Appendix E 87
 - 5.1.8 STEP 8: State the Result of Your Statistical Test 88
 - 5.1.9 STEP 9: State the Conclusion of Your Statistical Test in Plain English! 88
 - 5.2 Formula #1: Both Groups Have More Than 30 People in Them 92
 - 5.2.1 An Example of Formula #1 for the Two-Group t-Test 93
 - 5.3 Formula #2: One or Both Groups Have Less Than 30 People in Them 99
 - 5.4 End-of-Chapter Practice Problems 105
 - References 107
- 6 Correlation and Simple Linear Regression 109**
 - 6.1 What Is a “Correlation?” 109
 - 6.1.1 Understanding the Formula for Computing a Correlation 114
 - 6.1.2 Understanding the Nine Steps for Computing a Correlation, r 114

- 6.2 Using Excel to Compute a Correlation Between Two Variables 116
- 6.3 Creating a Chart and Drawing the Regression Line onto the Chart 121
 - 6.3.1 Using Excel to Create a Chart and the Regression Line Through the Data Points 123
- 6.4 Printing a Spreadsheet So That the Table and Chart Fit onto One Page 132
- 6.5 Finding the Regression Equation 133
 - 6.5.1 Installing the Data Analysis ToolPak into Excel 134
 - 6.5.2 Using Excel to Find the SUMMARY OUTPUT of Regression 136
 - 6.5.3 Finding the Equation for the Regression Line 140
 - 6.5.4 Using the Regression Line to Predict the y-Value for a Given x-Value 141
- 6.6 Adding the Regression Equation to the Chart 142
- 6.7 How to Recognize Negative Correlations in the SUMMARY OUTPUT Table 145
- 6.8 Printing Only Part of a Spreadsheet Instead of the Entire Spreadsheet 145
 - 6.8.1 Printing Only the Table and the Chart on a Separate Page 146
 - 6.8.2 Printing Only the Chart on a Separate Page 146
 - 6.8.3 Printing Only the SUMMARY OUTPUT of the Regression Analysis on a Separate Page 147
- 6.9 End-of-Chapter Practice Problems 147
- References 152
- 7 Multiple Correlation and Multiple Regression 153**
 - 7.1 Multiple Regression Equation 153
 - 7.2 Finding the Multiple Correlation and the Multiple Regression Equation 156
 - 7.3 Using the Regression Equation to Predict Annual Sales 159
 - 7.4 Using Excel to Create a Correlation Matrix in Multiple Regression 160
 - 7.5 End-of-Chapter Practice Problems 163
 - References 168
- 8 One-Way Analysis of Variance (ANOVA) 169**
 - 8.1 Using Excel to Perform a One-Way Analysis of Variance (ANOVA) 171
 - 8.2 How to Interpret the ANOVA Table Correctly 173
 - 8.3 Using the Decision Rule for the ANOVA F-Test 174

8.4	Testing the Difference Between Two Groups Using the ANOVA t-Test	175
8.4.1	Comparing Dierberg’s vs. Shop ‘n Save in Their Prices Using the ANOVA t-Test	175
8.5	End-of-Chapter Practice Problems	180
	References	187
Appendices	189
	Appendix A: Answers to End-of-Chapter Practice Problems	189
	Appendix B: Practice Test	222
	Appendix C: Answers to Practice Test	234
	Appendix D: Statistical Formulas	244
	Appendix E: t-TABLE	246
Index	247

Chapter 1

Sample Size, Mean, Standard Deviation, and Standard Error of the Mean



This chapter deals with how you can use Excel to find the average (i.e., “mean”) of a set of scores, the standard deviation of these scores (STDEV), and the standard error of the mean (s.e.) of these scores. All three of these statistics are used frequently and form the basis for additional statistical tests.

1.1 Mean

The *mean* is the “arithmetic average” of a set of scores. When my daughter was in the fifth grade, she came home from school with a sad face and said that she didn’t get “averages.” The book she was using described how to find the mean of a set of scores, and so I said to her:

“Jennifer, you add up all the scores and divide by the number of numbers that you have.”
She gave me “that look,” and said: “Dad, this is serious!” She thought I was teasing her.
So I said:
“See these numbers in your book; add them up. What is the answer?” (She did that.)
“Now, how many numbers do you have?” (She answered that question.)
“Then, take the number you got when you added up the numbers, and divide that number by the number of numbers that you have.”

She did that, and found the correct answer. You will use that same reasoning now, but it will be much easier for you because Excel will do all of the steps for you.

We will call this average of the scores the “mean” which we will symbolize as: \bar{X} , and we will pronounce it as: “Xbar.”

The formula for finding the mean with your calculator looks like this:

$$\bar{X} = \frac{\sum X}{n} \tag{1.1}$$

The symbol Σ is the Greek letter sigma, which stands for “sum.” It tells you to add up all the scores that are indicated by the letter X , and then to divide your answer by n (the number of numbers that you have).

Let’s give a simple example:

Suppose that you had these six scores:

6
4
5
3
2
5

To find the mean of these scores, you add them up, and then divide by the number of scores. So, the mean is: $25/6 = 4.17$.

1.2 Standard Deviation

The *standard deviation* tells you “how close the scores are to the mean.” If the standard deviation is a small number, this tells you that the scores are “bunched together” close to the mean. If the standard deviation is a large number, this tells you that the scores are “spread out” a greater distance from the mean. The formula for the standard deviation (which we will call *STDEV*) and use the letter, S , to symbolize is:

$$STDEV = S = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n - 1}} \quad (1.2)$$

The formula look complicated, but what it asks you to do is this:

1. Subtract the mean from each score ($X - \bar{X}$).
2. Then, square the resulting number to make it a positive number.
3. Then, add up these squared numbers to get a total score.
4. Then, take this total score and divide it by $n - 1$ (where n stands for the number of numbers that you have).
5. The final step is to take the square root of the number you found in step 4.

You will not be asked to compute the standard deviation using your calculator in this book, but you could see examples of how it is computed in any basic statistics book. Instead, we will use Excel to find the standard deviation of a set of scores. When we use Excel on the six numbers we gave in the description of the mean above, you will find that the *STDEV* of these numbers, S , is 1.47.

1.3 Standard Error of the Mean

The formula for the *standard error of the mean* (*s.e.*, which we will use $S_{\bar{X}}$ to symbolize) is:

$$\text{s.e.} = S_{\bar{X}} = \frac{S}{\sqrt{n}} \quad (1.3)$$

To find *s.e.*, all you need to do is to take the standard deviation, STDEV, and divide it by the square root of n, where n stands for the “number of numbers” that you have in your data set. In the example under the standard deviation description above, the *s.e.* = 0.60. (You can check this on your calculator.)

If you want to learn more about the standard deviation and the standard error of the mean, see Weiers (2011).

Now, let’s learn how to use Excel to find the sample size, the mean, the standard deviation, and the standard error of the mean using a problem from sales:

Suppose that you wanted to estimate the first-year sales of a new product that your company was about to launch into the marketplace. You have decided to look at the first-year sales of similar products that your company has launched to get an idea of what sales are typical for your new product launches.

You decide to use the first-year sales of a similar product over the past 8 years, and you have created the table in Fig. 1.1:

Fig. 1.1 Worksheet Data for First-year Sales (Practical Example)

Year	First-year sales (\$000)
1	10
2	10
3	12
4	16
5	22
6	29
7	39
8	47

Note that the first-year sales are in thousands of dollars (\$000), so that 10 means that the first-year sales of that product were really \$10,000.

1.4 Sample Size, Mean, Standard Deviation, and Standard Error of the Mean

Objective: To find the sample size (n), mean, standard deviation (STDEV), and standard error of the mean (s.e.) for these data

Start your computer, and click on the Excel 2019 icon to open a blank Excel spreadsheet.

Click on: Blank workbook

Enter the data in this way:

A3: Year

B3: First-year sales (\$000)

A4: 1

1.4.1 Using the Fill/Series/Columns Commands

Objective: To add the years 2–8 in a column underneath year 1

Put pointer in A4

Home (top left of screen)

Important note: The “Paste” icon should be on the top of your screen on the far left of the screen.

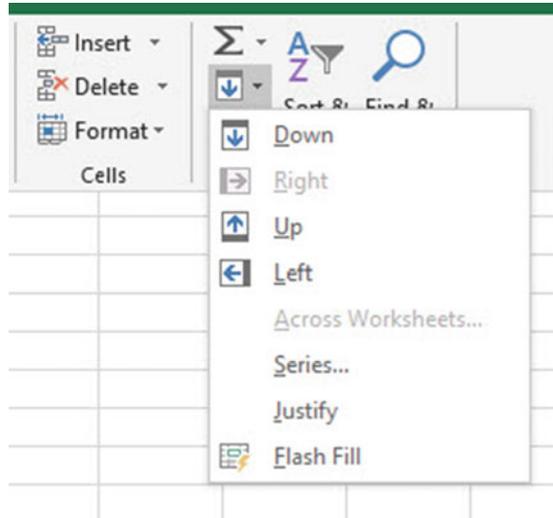
Important note: Notice the Excel commands at the top of your computer screen:

File→**Home**→**Insert**→**Page Layout**→**Formulas**→*etc.*

If these commands ever “disappear” when you are using Excel, you need to click on “Home” at the top of your screen to make them reappear!

Fill (top right of screen: click on the down arrow; see Fig. 1.2)

Fig. 1.2 Home/Fill/Series commands



Series
 Columns
 Step value: 1
 Stop value: 8 (see Fig. 1.3)

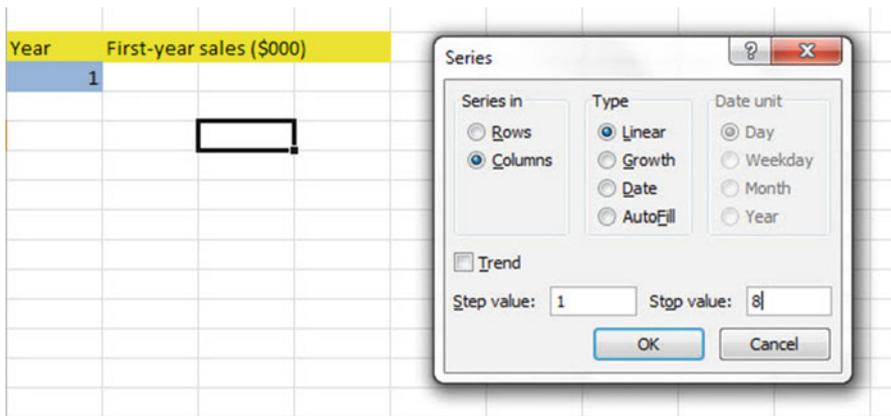


Fig. 1.3 Example of Dialog Box for Fill/Series/Columns/Step Value/Stop Value commands

OK

The years should be identified as 1–8, with 8 in cell A11.

Now, enter the first-year sales figures in cells B4: B11 using the above table.

Since your computer screen shows the information in a format that does not look professional, you need to learn how to “widen the column width” and how to “center the information” in a group of cells. Here is how you can do those two steps:

1.4.2 Changing the Width of a Column

Objective: To make a column width wider so that all of the information fits inside that column

If you look at your computer screen, you can see that Column B is not wide enough so that all of the information fits inside this column. To make Column B wider:

Click on the letter, B, at the top of your computer screen

Place your mouse pointer at the far right corner of B until you create a “cross sign” on that corner

Left-click on your mouse, hold it down, and move this corner to the right until it is “wide enough to fit all of the data”

Take your finger off the mouse to set the new column width (see Fig. 1.4)

Fig. 1.4 Example of How to Widen the Column Width

A	B	C
Year	First-year sales (\$000)	
1		10
2		10
3		12
4		16
5		22
6		29
7		39
8		47

Then, click on any empty cell (i.e., any blank cell) to “deselect” column B so that it is no longer a darker color on your screen.

When you widen a column, you will make all of the cells in all of the rows of this column that same width.

Now, let’s go through the steps to center the information in both Column A and Column B.

1.4.3 Centering Information in a Range of Cells

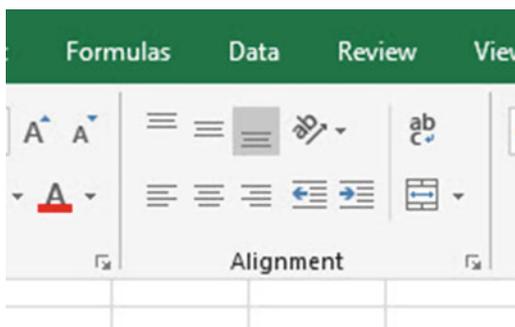
Objective: To center the information in a group of cells

In order to make the information in the cells look “more professional,” you can center the information using the following steps:

Left-click your mouse on A3 and drag it to the right and down to highlight cells A3: B11 so that these cells appear in a darker color
Home

At the top of your computer screen, you will see a set of “lines” in which all of the lines are “centered” to the same width under “Alignment” (it is the second icon at the bottom left of the Alignment box; see Fig. 1.5)

Fig. 1.5 Example of How to Center Information Within Cells



Click on this icon to center the information in the selected cells (see Fig. 1.6)

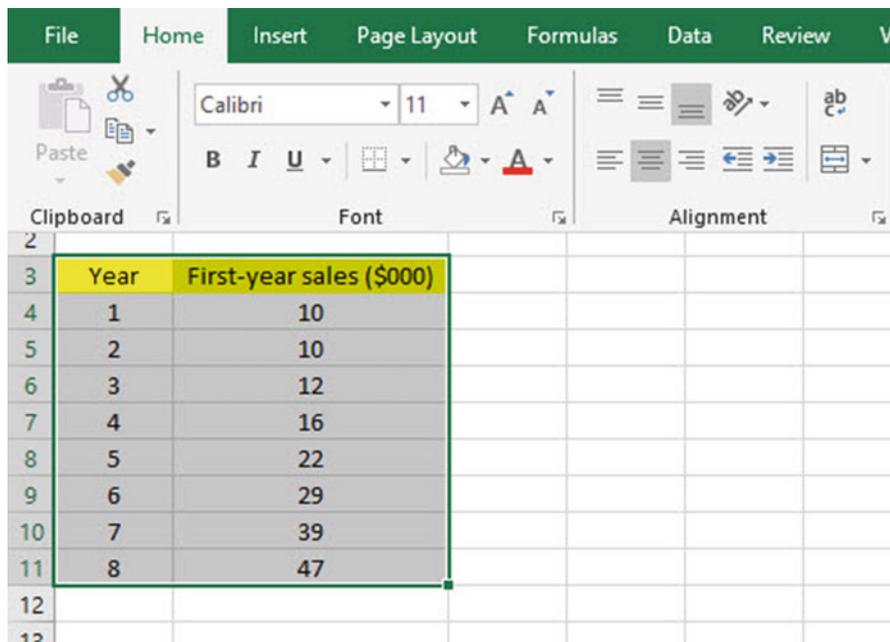


Fig. 1.6 Final Result of Centering Information in the Cells

Since you will need to refer to the first-year sales figures in your formulas, it will be much easier to do this if you “name the range of data” with a name instead of having to remember the exact cells (B4: B11) in which these figures are located. Let’s call that group of cells Product, but we could give them any name that you want to use.

1.4.4 Naming a Range of Cells

Objective: To name the range of data for the first-year sales figures with the name: Product

Highlight cells B4: B11 by left-clicking your mouse on B4 and dragging it down to B11

Formulas (top left of your screen)

Define Name (top center of your screen)

Product (type this name in the top box; see Fig. 1.7)

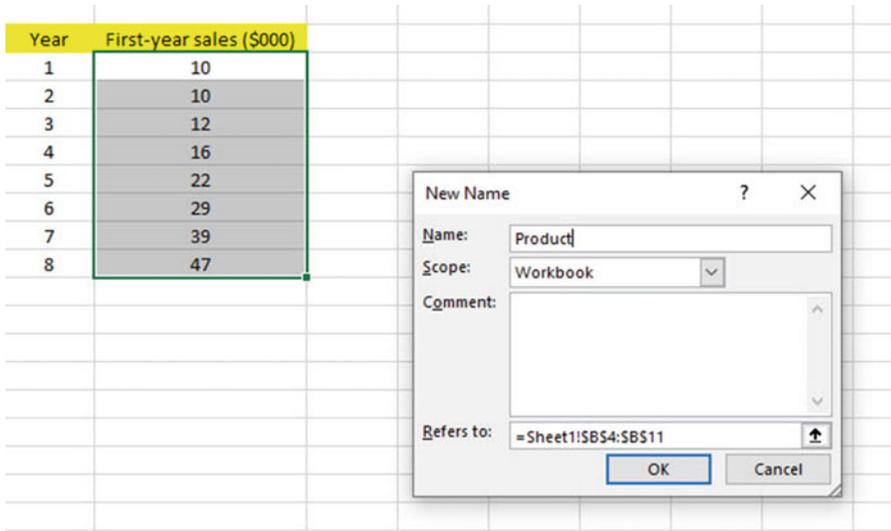


Fig. 1.7 Dialogue box for “naming a range of cells” with the name: Product

OK

Then, click on any cell of your spreadsheet that does not have any information in it (i.e., it is an “empty cell”) to deselect cells B4:B11

Now, add the following terms to your spreadsheet:

- E6: n
- E9: Mean
- E12: STDEV
- E15: s.e. (see Fig. 1.8)

	A	B	C	D	E	F
1						
2						
3	Year	First-year sales (\$000)				
4	1	10				
5	2	10				
6	3	12			n	
7	4	16				
8	5	22				
9	6	29			Mean	
10	7	39				
11	8	47				
12					STDEV	
13						
14						
15					s.e.	
16						

Fig. 1.8 Example of Entering the Sample Size, Mean, STDEV, and s.e. Labels

Note: Whenever you use a formula, you must add an equal sign (=) at the beginning of the name of the function so that Excel knows that you intend to use a formula.

1.4.5 Finding the Sample Size Using the =COUNT Function

Objective: To find the sample size (n) for these data using the =COUNT function

F6: =COUNT(Product)

Hit the Enter key, and this command should insert the number 8 into cell F6 since there are eight first-year sales figures.

1.4.6 Finding the Mean Score Using the =AVERAGE Function

Objective: To find the mean sales figure using the =AVERAGE function

F9: =AVERAGE(Product)

This command should insert the number 23.125 into cell F9.

1.4.7 Finding the Standard Deviation Using the =STDEV Function

Objective: To find the standard deviation (STDEV) using the =STDEV function

F12: =STDEV(Product)

This command should insert the number 14.02485 into cell F12.

1.4.8 Finding the Standard Error of the Mean

Objective: To find the standard error of the mean using a formula for these eight data points

F15: =F12/SQRT(8)

This command should insert the number 4.958533 into cell F15 (see Fig. 1.9).

	A	B	C	D	E	F	G
1							
2							
3	Year	First-year sales (\$000)					
4	1	10					
5	2	10					
6	3	12			n	8	
7	4	16					
8	5	22					
9	6	29			Mean	23.125	
10	7	39					
11	8	47					
12					STDEV	14.02485	
13							
14							
15					s.e.	4.958533	
16							

Fig. 1.9 Example of Using Excel Formulas for Sample Size, Mean, STDEV, and s.e.

Important note: Throughout this book, be sure to double-check all of the figures in your spreadsheet to make sure that they are in the correct cells, or the formulas will not work correctly!

1.4.8.1 Formatting Numbers in Number Format (Two Decimal Places)

Objective: To convert the mean, STDEV, and s.e. to two decimal places

Highlight cells F9:F15

Click on: Home (top left of screen)

Click on the down arrow to the right of “Number” at the top center of your screen.

Inside the dialog box, click on; Number

Keep the two decimal places already selected (see Fig. 1.10)

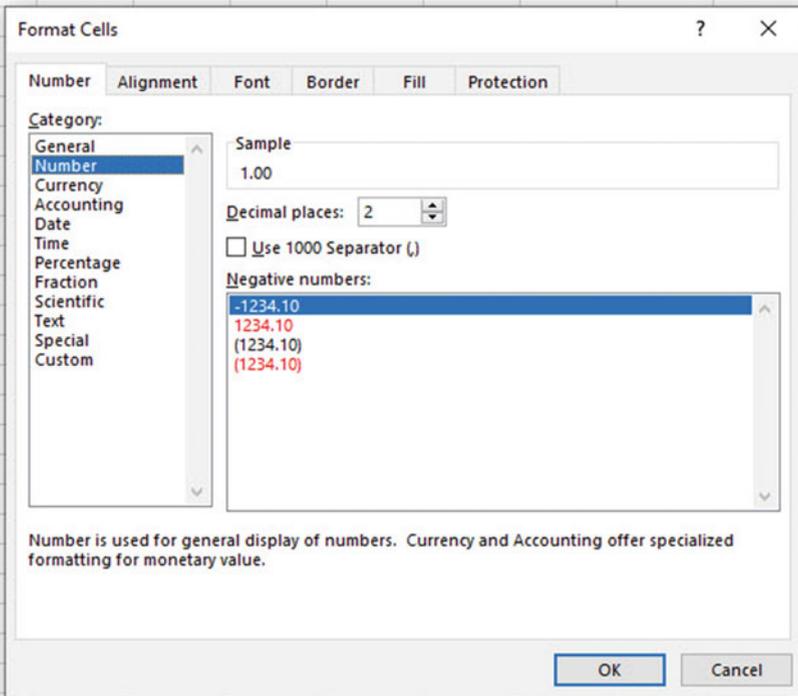


Fig. 1.10 Using the Number Format dialog box to convert Numbers to two Decimal Places

OK (see Fig. 1.11)

Year	First-year sales (\$000)		
1	10		
2	10		
3	12	n	8
4	16		
5	22		
6	29	Mean	23.13
7	39		
8	47	STDEV	14.02
		s.e.	4.96

Fig. 1.11 Example of Converting Numbers to Two Decimal Places

Important note: The sales figures are in thousands of dollars (\$000), so that the mean is \$23,130, the standard deviation is \$14,020, and the standard error of the mean is \$4960.

Now, click on any “empty cell” on your spreadsheet to deselect cells F9:F15.

1.5 Saving a Spreadsheet

Objective: To save this spreadsheet with the name: Product6

In order to save your spreadsheet so that you can retrieve it sometime in the future, your first decision is to decide “where” you want to save it. That is your decision and you have several choices. If it is your own computer, you can save it onto your hard drive (you need to ask someone how to do that on your computer). Or, you can save it onto a “CD” or onto a “flash drive.”

When you want to save your file into your “Documents” location on your computer, you then need to complete these steps:

Click on: File (top of screen, far left)

Save as

This PC

Click on: Documents

(These commands will save this file in: “This PC: Documents location”)

File name: Product6 (enter this name to the right of File name; see Fig. 1.12)

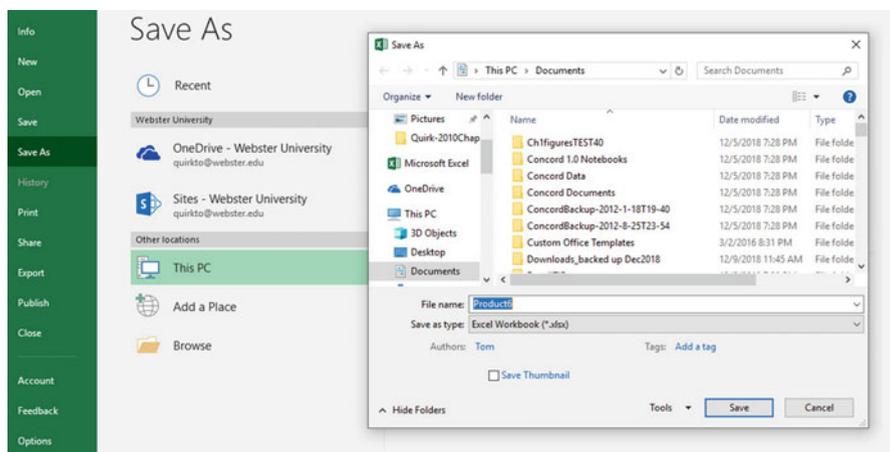


Fig. 1.12 Dialogue Box of Saving an Excel Workbook File as “Product6” in Documents location

Save (bottom right of dialog box)

Important note: Be very careful to save your Excel file spreadsheet every few minutes so that you do not lose your information!

1.6 Printing a Spreadsheet

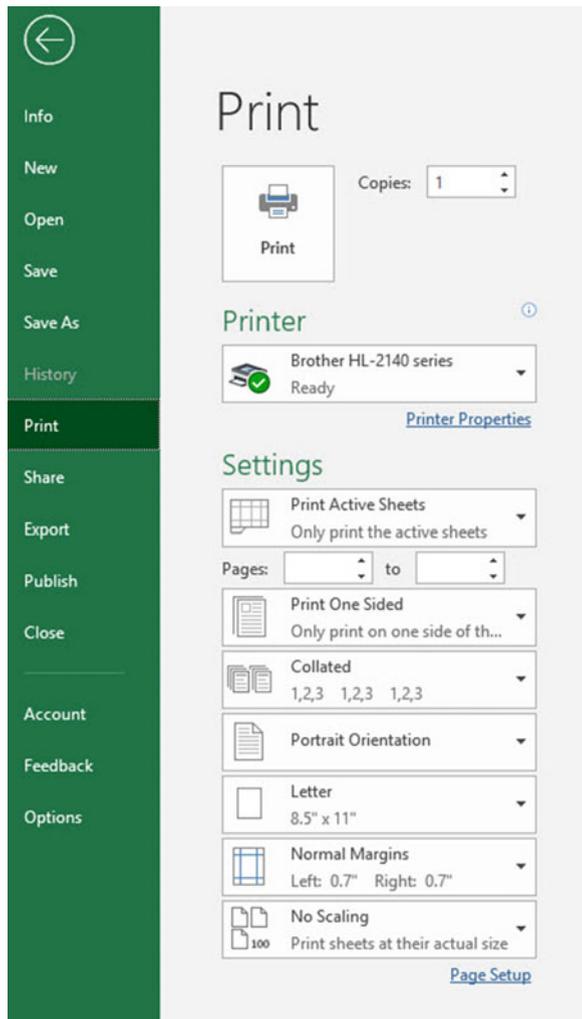
Objective: To print the spreadsheet

Use the following procedure when printing any spreadsheet.

File (top of screen, far left icon)

Print (see Fig. 1.13)

Fig. 1.13 Example of How to Print an Excel Worksheet Using the File/Print Commands



Print (at top left of screen)

The final spreadsheet is given in Fig. 1.14

Year	First-year sales (\$000)			
1	10			
2	10			
3	12		n	8
4	16			
5	22			
6	29		Mean	23.13
7	39			
8	47			
			STDEV	14.02
			s.e.	4.96

Fig. 1.14 Final Result of Printing an Excel Spreadsheet

Before you leave this chapter, let’s practice changing the format of the figures on a spreadsheet with two examples: (1) using two decimal places for figures that are dollar amounts, and (2) using three decimal places for figures.

Save the final spreadsheet by: File/Save, then close your spreadsheet by: File/Close, and then open a blank Excel spreadsheet by using: File/New/Blank Workbook icon (on the top left of your screen).

1.7 Formatting Numbers in Currency Format (Two Decimal Places)

Objective: To change the format of figures to dollar format with two decimal places

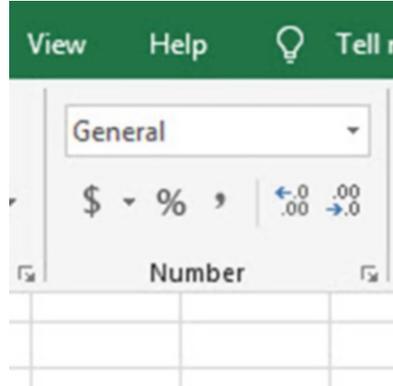
- A3: Price
- A4: 1.25
- A5: 3.45
- A6: 12.95

Highlight cells A4:A6 by left-clicking your mouse on A4 and dragging it down so that these three cells are highlighted in a darker color

Home

Number (top center of screen: click on the down arrow on the right; see Fig. 1.15)

Fig. 1.15 Dialogue Box for Number Format Choices



Category: Currency
Decimal places: 2 (then see Fig. 1.16)

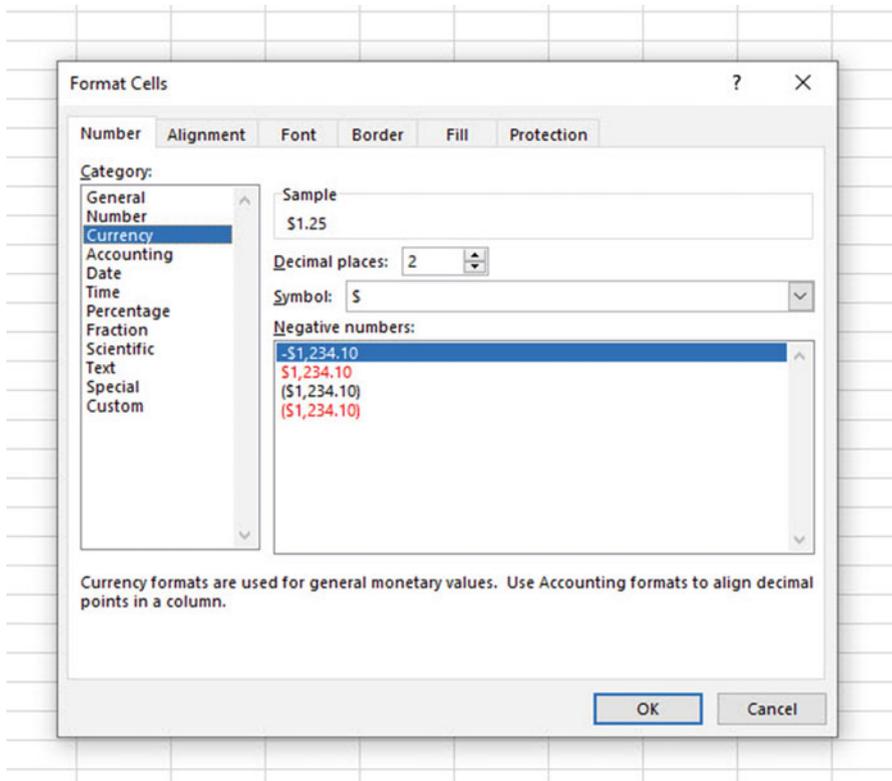


Fig. 1.16 Dialogue Box for Currency (two decimal places) Format for Numbers

OK

The three cells should have a dollar sign in them and be in two decimal places. Next, let's practice formatting figures in number format, three decimal places.

1.8 Formatting Numbers in Number Format (Three Decimal Places)

Objective: To format figures in number format, three decimal places

Home

Highlight cells A4:A6 on your computer screen

Number (click on the down arrow on the right)

Category: number

At the right of the box, change two decimal places to three decimal places by clicking on the "up arrow" once

OK

The three figures should now be in number format, each with three decimals.

Now, click on any blank cell to deselect cells A4:A6. Then, close this file by File/Close/Don't Save (since there is no need to save this practice problem).

You can use these same commands to format a range of cells in percentage format (and many other formats) to whatever number of decimal places you want to specify.

1.9 End-of-Chapter Practice Problems

1. Suppose that you work for Ford Motor Company and that you have been asked to do the data analysis for a panel of female college students to determine their importance of possible features for a new Ford Focus that is on the drawing board. You want to test your Excel skills, and so you have created a table of hypothetical data for Item #12 in the Survey that will be administered to this panel over the Web. These data are given in Fig. 1.17:

Survey of new-car features						
Panel of female college students (ages 18-24)						
Question #12: If you were to purchase a new car today, how important to you is the feature that "the car parallel parks itself to the curb" by using a computer?						
1	2	3	4	5	6	7
Not Important						Very Important
		RATING				
		5				
		6				
		4				
		3				
		7				
		6				
		5				
		7				
		6				
		7				
		4				
		3				
		1				
		7				
		6				
		4				
		5				

Fig. 1.17 Worksheet Data for Chap. 1: Practice Problem #1

- (a) Use Excel to the right of the table to find the sample size, mean, standard deviation, and standard error of the mean for these data. Label your answers, and round off the mean, standard deviation, and standard error of the mean to three decimal places; use number format for these three figures.
 - (b) Print the result on a separate page.
 - (c) Save the file as: CAR12A
2. Suppose that the Human Resources department of your company has administered a "Morale Survey" to all middle-level managers and that you have been asked to summarize the results of the survey. You have decided to test your Excel skills on one item to see if you can do this assignment correctly, and you have selected item #21 to test out your skills. The hypothetical data are given in Fig. 1.18.

HUMAN RESOURCES MORALE SURVEY						
Item #21: "Management is doing a good job of keeping employee morale at a high level."						
1	2	3	4	5	6	7
Disagree						Agree
			Rating			
			3			
			6			
			5			
			7			
			2			
			3			
			6			
			5			
			4			
			7			
			6			
			1			
			3			
			2			
			4			
			5			
			6			
			4			
			5			
			3			
			6			
			4			
			7			

Fig. 1.18 Worksheet Data for Chap. 1: Practice Problem #2

- (a) Use Excel to create a table of these ratings, and at the right of the table use Excel to find the sample size, mean, standard deviation, and standard error of the mean for these data. Label your answers, and round off the mean, standard deviation, and standard error of the mean to two decimal places using number format.
 - (b) Print the result on a separate page.
 - (c) Save the file as: MORALE4
3. Suppose that you have been hired to do analysis of data from the previous 18 days at a Ford assembly plant that produces Ford Focus automobiles. The plant manager wants you to summarize the number of defects per day of this car

produced during this 3-week period. A “defect” is defined as any irregularity of the car at the end of the production line that requires the car to be brought off the line and repaired before it is shipped to a dealer. The data from the previous 3 weeks are given in Fig. 1.19:

Ford Motor Co.	
Number of defects per day for the Ford Focus	
Day	No. of defects
1	6
2	8
3	14
4	12
5	6
6	8
7	23
8	17
9	14
10	16
11	18
12	12
13	13
14	15
15	8
16	6
17	9
18	10

Fig. 1.19 Worksheet Data for Chap. 1: Practice Problem #3

- Use Excel to create a table for these data, and at the right of the table, use Excel to find the sample size, mean, standard deviation, and standard error of the mean for these data. Label your answers, and round off the mean, standard deviation, and standard error of the mean to three decimal places using number format.
- Print the result on a separate page.
- Save the file as: DEFECTS4

Reference

Weiers, R.M. Introduction to Business Statistics (7th ed.). Mason, OH: South-Western Cengage Learning, 2011.

Chapter 2

Random Number Generator



Suppose that you wanted to take a random sample of 5 of your company’s 32 salespeople using Excel so that you could interview these five salespeople about their job satisfaction at your company.

To do that, you need to define a “sampling frame.” A sampling frame is a list of people from which you want to select a random sample. This frame starts with the identification code (ID) of the number 1 that is assigned to the name of the first salesperson in your list of 32 salespeople in your company. The second salesperson has a code number of 2, the third a code number of 3, and so forth until the last salesperson has a code number of 32.

Since your company has 32 salespeople, your sampling frame would go from 1 to 32 with each salesperson having a unique ID number.

We will first create the frame numbers as follows in a new Excel worksheet:

2.1 Creating Frame Numbers for Generating Random Numbers

Objective: To create the frame numbers for generating random numbers

A3: FRAME NO.

A4: 1

Now, create the frame numbers in column A with the Home/Fill commands that were explained in the first chapter of this book (see Sect. 1.4.1) so that the frame numbers go from 1 to 32, with the number 32 in cell A35. If you need to be reminded about how to do that, here are the steps:

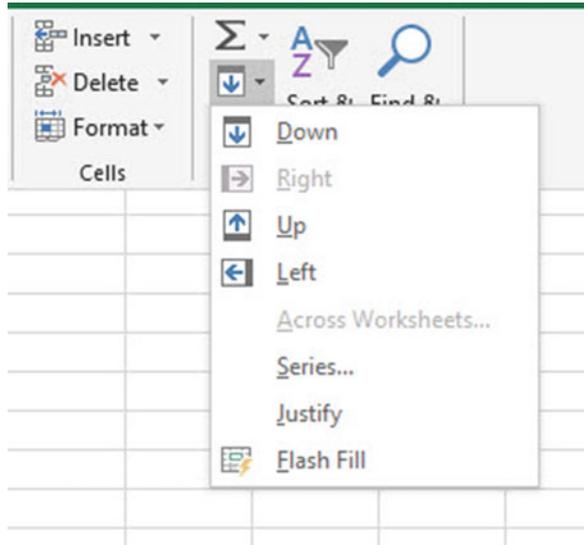
Click on cell A4 to select this cell

Home

Fill (then click on the “down arrow” next to this command and select)

Series (see Fig. 2.1)

Fig. 2.1 Dialogue Box for Fill/Series Commands



Columns

Step value: 1

Stop value: 32 (see Fig. 2.2)

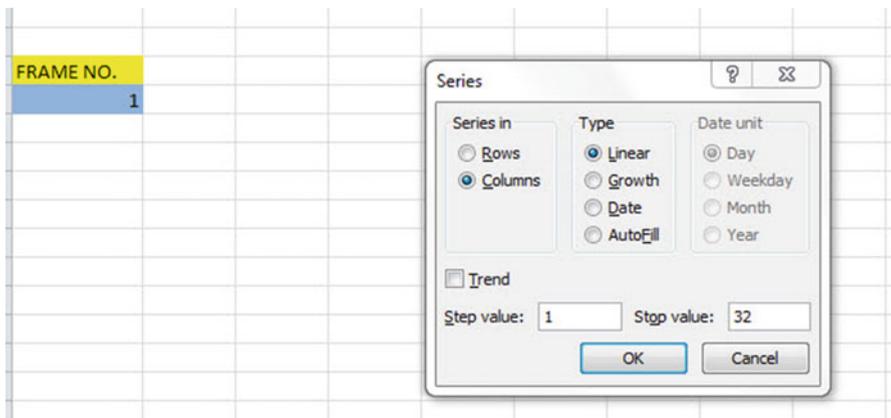


Fig. 2.2 Dialogue Box for Fill/Series/Columns/Step value/Stop value Commands

OK

Then, save this file as: Random2. You should obtain the result in Fig. 2.3.

Fig. 2.3 Frame Numbers from 1 to 32

FRAME NO.			
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			
26			
27			
28			
29			
30			
31			
32			

Now, create a column next to these frame numbers in this manner:

B3: DUPLICATE FRAME NO.

B4: 1

Next, use the Home/Fill command again, so that the 32 frame numbers begin in cell B4 and end in cell B35. Be sure to widen the columns A and B so that all of the information in these columns fits inside the column width. Then, center the information inside both Column A and Column B on your spreadsheet. You should obtain the information given in Fig. 2.4.

FRAME NO.	DUPLICATE FRAME NO.			
1	1			
2	2			
3	3			
4	4			
5	5			
6	6			
7	7			
8	8			
9	9			
10	10			
11	11			
12	12			
13	13			
14	14			
15	15			
16	16			
17	17			
18	18			
19	19			
20	20			
21	21			
22	22			
23	23			
24	24			
25	25			
26	26			
27	27			
28	28			
29	29			
30	30			
31	31			
32	32			

Fig. 2.4 Duplicate Frame Numbers from 1 to 32

Save this file as: Random3

You are probably wondering why you created the same information in both Column A and Column B of your spreadsheet. This is to make sure that before you sort the frame numbers that you have exactly 32 of them when you finish sorting them into a random sequence of 32 numbers.

Now, let's add a random number to each of the duplicate frame numbers as follows:

2.2 Creating Random Numbers in an Excel Worksheet

C3: RANDOM NO.

(then widen columns A, B, C so that their labels fit inside the columns; then center the information in A3:C35)

C4: =RAND()

Next, hit the Enter key to add a random number to cell C4.

Note that you need *both* an open parenthesis *and* a closed parenthesis after =RAND(). The RAND command “looks to the left of the cell with the RAND() COMMAND in it” and assigns a random number to that cell.

Now, put the pointer using your mouse in cell C4 and then move the pointer to the bottom right corner of that cell until you see a “plus sign” in that cell. Then, click and drag the pointer down to cell C35 to add a random number to all 32 ID frame numbers (see Fig. 2.5).

FRAME NO.	DUPLICATE FRAME NO.	RANDOM NO.		
1	1	0.690332931		
2	2	0.022334603		
3	3	0.89452184		
4	4	0.981573849		
5	5	0.698381228		
6	6	0.611413628		
7	7	0.013551391		
8	8	0.036862479		
9	9	0.412932328		
10	10	0.460808373		
11	11	0.533416136		
12	12	0.988470378		
13	13	0.097821358		
14	14	0.881481661		
15	15	0.352287507		
16	16	0.344014139		
17	17	0.084570168		
18	18	0.467909507		
19	19	0.904917153		
20	20	0.252482436		
21	21	0.788783634		
22	22	0.592964999		
23	23	0.946665187		
24	24	0.214249616		
25	25	0.509340791		
26	26	0.439105519		
27	27	0.086378662		
28	28	0.975489923		
29	29	0.120077924		
30	30	0.216062043		
31	31	0.353995884		
32	32	0.558171248		

Fig. 2.5 Example of Random Numbers Assigned to the Duplicate Frame Numbers

Then, click on any empty cell to deselect C4:C35 to remove the dark color highlighting these cells.

Save this file as: Random3A

Now, let's sort these duplicate frame numbers into a random sequence:

2.3 Sorting Frame Numbers into a Random Sequence

Objective: To sort the duplicate frame numbers into a random sequence

Highlight cells B3: C35 (include the labels at the top of columns B and C)
Data (top of screen)
Sort (click on this word at the top center of your screen; see Fig. 2.6)

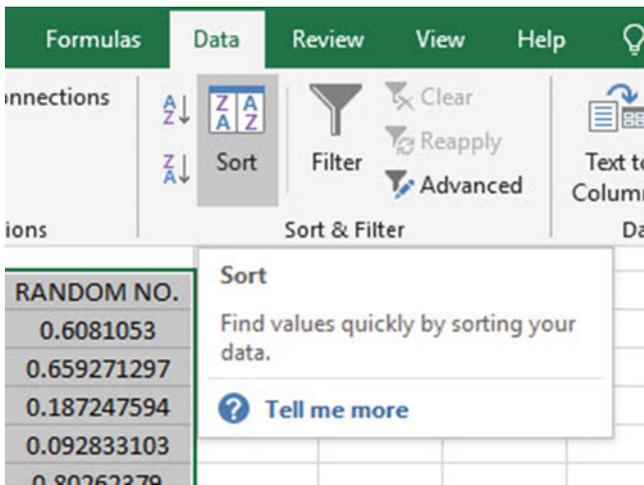


Fig. 2.6 Dialogue Box for Data/Sort Commands

Sort by: RANDOM NO. (click on the down arrow)
Smallest to Largest (see Fig. 2.7)

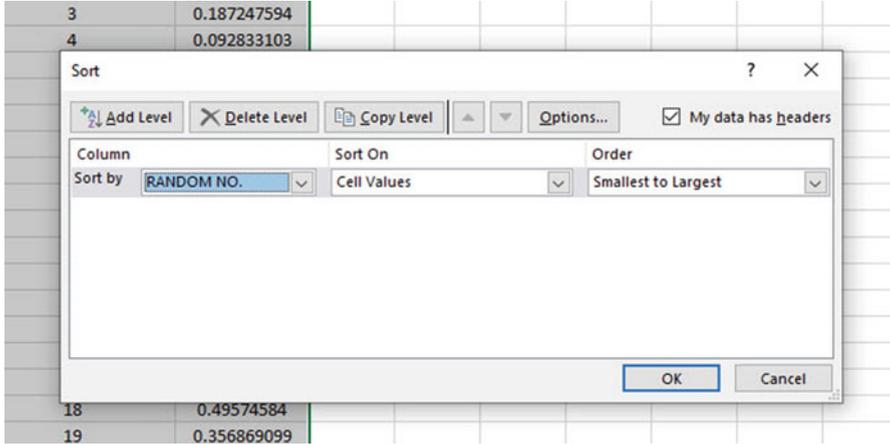


Fig. 2.7 Dialogue Box for Data/Sort/RANDOM NO./Smallest to Largest Commands

OK

Click on any empty cell to deselect B3:C35.

Save this files as: Random4

Print this file now.

These steps will produce Fig. 2.8 with the DUPLICATE FRAME NUMBER sorted into a random order:

FRAME NO.	DUPLICATE FRAME NO.	RANDOM NO.
1	7	0.343261283
2	2	0.929607291
3	8	0.914304212
4	17	0.903618324
5	27	0.257228182
6	13	0.456204036
7	29	0.390622986
8	24	0.222210116
9	30	0.432155483
10	20	0.219982266
11	16	0.842461398
12	15	0.3781508
13	31	0.694049089
14	9	0.939764564
15	26	0.075689667
16	10	0.302227714
17	18	0.468687794
18	25	0.148502036
19	11	0.49462371
20	32	0.87719372
21	22	0.413151766
22	6	0.094310793
23	1	0.962115342
24	5	0.528964967
25	21	0.401140496
26	14	0.403327013
27	3	0.865025638
28	19	0.517332393
29	23	0.968085821
30	28	0.647609375
31	4	0.670143403
32	12	0.09483352

Fig. 2.8 Duplicate Frame Numbers Sorted by Random Number

Important note: Because Excel randomly assigns these random numbers, your Excel commands will produce a different sequence of random numbers from everyone else who reads this book!

Because your objective at the beginning of this chapter was to select randomly 5 of your company’s 32 salespeople for a personal interview, you now can do that by selecting the *first five ID numbers* in DUPLICATE FRAME NO. column after the sort.

Although your first five random numbers will be different from those we have selected in the random sort that we did in this chapter, we would select these five IDs of salespeople to interview using Fig. 2.9.

7, 2, 8, 17, 27

FRAME NO.	DUPLICATE FRAME NO.	RANDOM NO.
1	7	0.343261283
2	2	0.929607291
3	8	0.914304212
4	17	0.903618324
5	27	0.257228182
6	13	0.456204036
7	29	0.390622986
8	24	0.222210116
9	30	0.432155483
10	20	0.219982266
11	16	0.842461398
12	15	0.3781508
13	31	0.694049089
14	9	0.939764564
15	26	0.075689667
16	10	0.302227714
17	18	0.468687794
18	25	0.148502036
19	11	0.49462371
20	32	0.87719372
21	22	0.413151766
22	6	0.094310793
23	1	0.962115342
24	5	0.528964967
25	21	0.401140496
26	14	0.403327013
27	3	0.865025638
28	19	0.517332393
29	23	0.968085821
30	28	0.647609375
31	4	0.670143403
32	12	0.09483352

Fig. 2.9 First Five Salespeople Selected Randomly

Remember, your five ID numbers selected after your random sort will be different from the five ID numbers in Fig. 2.9 because Excel assigns a different random number *each time the =RAND() command is given*.

Before we leave this chapter, you need to learn how to print a file so that all of the information on that file fits onto a single page without “dribbling over” onto a second or third page.

2.4 Printing an Excel File So That All of the Information Fits onto One Page

Objective: To print a file so that all of the information fits onto one page

Note that the three practice problems at the end of this chapter require you to sort random numbers when the files contain 63 customers, 114 employees, and 76 key accounts, respectively. These files will be “too big” to fit onto one page when you print them unless you format these files so that they fit onto a single page when you print them.

Let’s create a situation where the file does not fit onto one printed page unless you format it first to do that.

Go back to the file you just created, Random 4, and enter the name: *Jennifer* into cell: A51.

If you printed this file now, the name, *Jennifer*, would be printed onto a second page because it “dribbles over” outside of the page range for this file in its current format.

So, you would need to change the page format so that all of the information, including the name, Jennifer, fits onto just one page when you print this file by using the following steps:

Page Layout (top left of the computer screen)

(Notice the “Scale to Fit” section in the center of your screen; see Fig. 2.10)

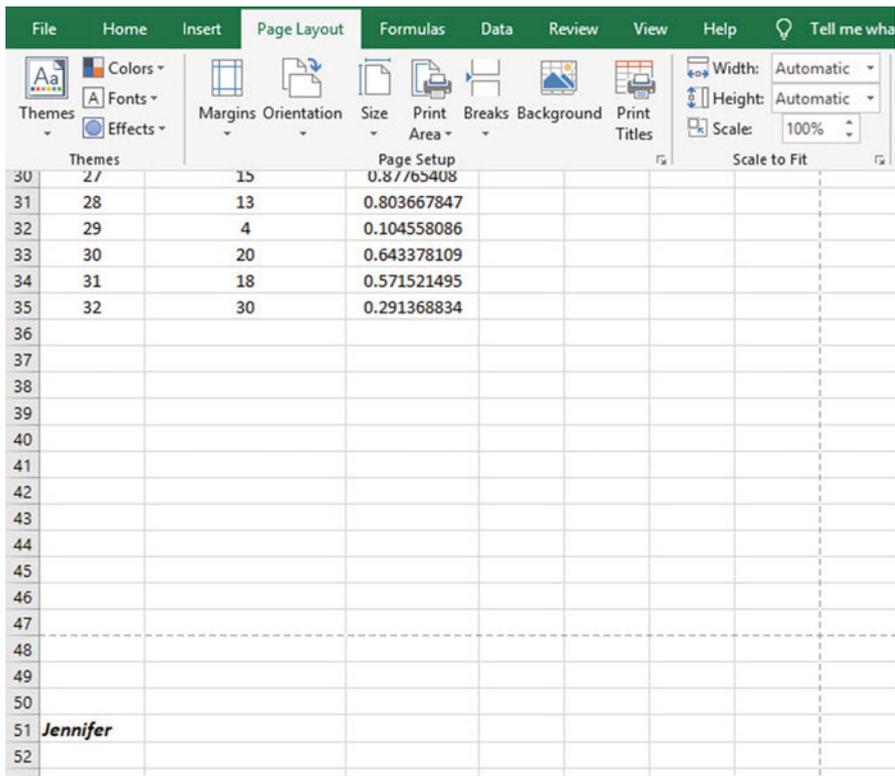


Fig. 2.10 Dialogue Box for Page Layout/Scale to Fit Commands

Hit the down arrow to the right of 100% *once* to reduce the size of the page to 95%

Now, note that the name, Jennifer, is still on a second page on your screen because her name is below the horizontal dotted line on your screen in Fig. 2.11 (the dotted lines tell you outline dimensions of the file if you printed it now).

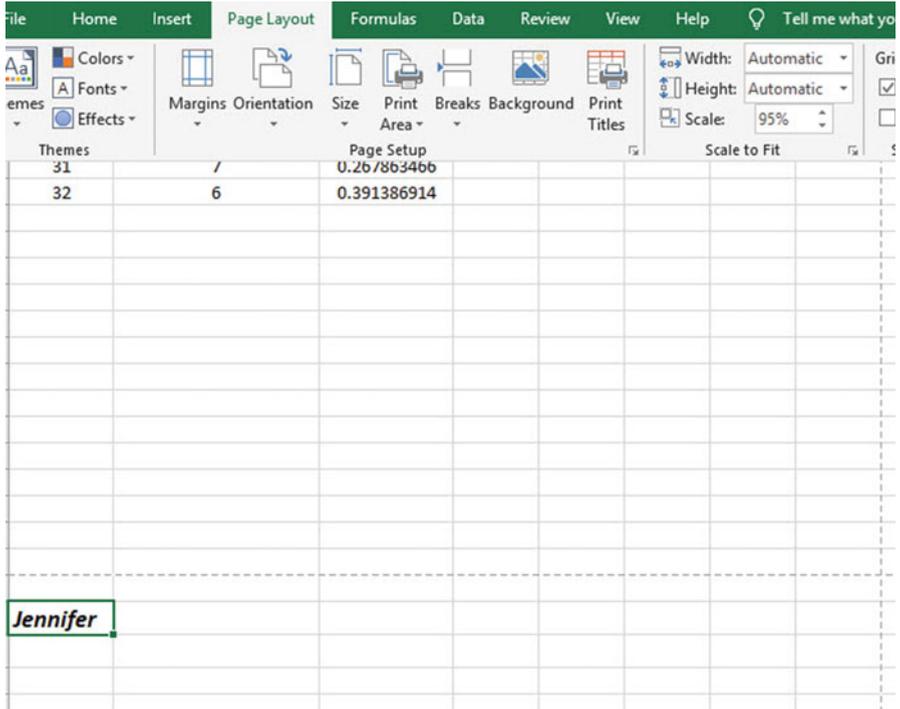


Fig. 2.11 Example of Scale Reduced to 95% with “Jennifer” to be Printed on a Second Page

So, you need to repeat the “scale change steps” by hitting the down arrow on the right once more to reduce the size of the worksheet to 90% of its normal size.

Notice that the “dotted lines” on your computer screen in Fig. 2.12 are now below Jennifer’s name to indicate that all of the information, including her name, is now formatted to fit onto just one page when you print this file.

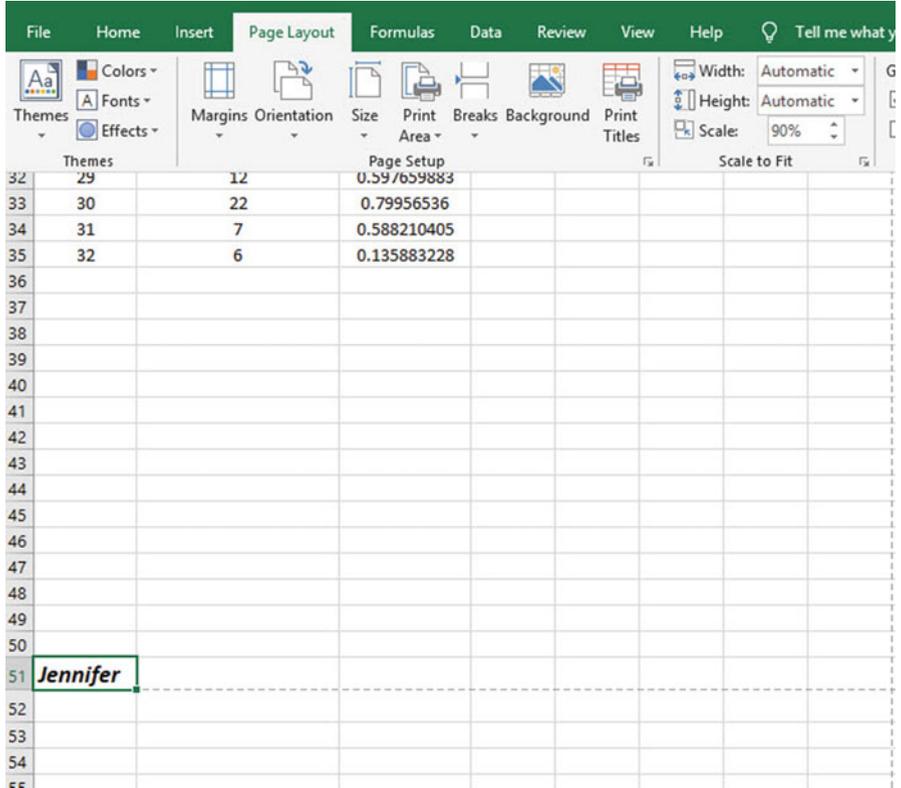


Fig. 2.12 Example of Scale Reduced to 90% with “Jennifer” to be printed on the first page (note the dotted line below Jennifer on your screen)

Save the file as: Random4A

Print the file. Does it all fit onto one page? It should (see Fig. 2.13).

2.5 End-of-Chapter Practice Problems

1. Suppose that you wanted to do a “customer satisfaction phone survey” of 15 of 63 customers who purchased at least \$1000 worth of merchandise from your company during the last 60 days.
 - (a) Set up a spreadsheet of frame numbers for these customers with the heading: FRAME NUMBERS using the Home/Fill commands.
 - (b) Then, create a separate column to the right of these frame numbers which duplicates these frame numbers with the title: Duplicate frame numbers
 - (c) Then, create a separate column to the right of these duplicate frame numbers and use the =RAND() function to assign random numbers to all of the frame numbers in the duplicate frame numbers column, and change this column format so that three decimal places appear for each random number
 - (d) Sort the duplicate frame numbers and random numbers into a random order
 - (e) Print the result so that the spreadsheet fits onto one page
 - (f) Circle on your printout the I.D. number of the first 15 customers that you would call in your phone survey
 - (g) Save the file as: RAND9

Important note: Note that everyone who does this problem will generate a different random order of customer ID numbers since Excel assigns a different random number each time the RAND() command is used. For this reason, the answer to this problem given in this Excel Guide will have a completely different sequence of random numbers from the random sequence that you generate. This is normal and what is to be expected.

2. Suppose that you are the Human Resources (HR) Director at a large company and that you want to conduct a phone interview with 10 of the 114 employees who elected to be included in the company’s Vision Care Plan which was instituted a year ago to obtain their feedback about how well the plan has been working for them.
 - (a) Set up a spreadsheet of frame numbers for these employees with the heading: FRAME NO.
 - (b) Then, create a separate column to the right of these frame numbers which duplicates these frame numbers with the title: Duplicate frame no.
 - (c) Then, create a separate column to the right of these duplicate frame numbers entitled “Random number” and use the =RAND() function to assign random numbers to all of the frame numbers in the duplicate frame numbers column. Then, change this column format so that three decimal places appear for each random number
 - (d) Sort the duplicate frame numbers and random numbers into a random order

- (e) Print the result so that the spreadsheet fits onto one page
 - (f) Circle on your printout the I.D. number of the first 10 employees that the HR Director would call in his phone survey
 - (g) Save the file as: RANDOM6
3. Suppose that your Sales department at your company wants to do a “customer satisfaction survey” of 20 of your company’s 76 “key accounts.” Suppose, further, that your Sales Vice-President has defined a key account as a customer who purchased at least \$30,000 worth of merchandise from your company in the past 90 days.
- (a) Set up a spreadsheet of frame numbers for these customers with the heading: FRAME NUMBERS.
 - (b) Then, create a separate column to the right of these frame numbers which duplicates these frame numbers with the title: Duplicate frame numbers
 - (c) Then, create a separate column to the right of these duplicate frame numbers entitled “Random number” and use the =RAND() function to assign random numbers to all of the frame numbers in the duplicate frame numbers column. Then, change this column format so that three decimal places appear for each random number
 - (d) Sort the duplicate frame numbers and random numbers into a random order
 - (e) Print the result so that the spreadsheet fits onto one page
 - (f) Circle on your printout the I.D. number of the first 20 customers that your Sales Vice-President would call for his phone survey.
 - (g) Save the file as: RAND5

Chapter 3

Confidence Interval About the Mean Using the TINV Function and Hypothesis Testing



This chapter focuses on two ideas: (1) finding the 95% confidence interval about the mean, and (2) hypothesis testing.

Let's talk about the confidence interval first.

3.1 Confidence Interval About the Mean

In statistics, we are always interested in *estimating the population mean*. How do we do that?

3.1.1 How to Estimate the Population Mean

Objective: To estimate the population mean, μ

Remember that the population mean is the average of all of the people in the target population. For example, if we were interested in how well adults ages 25–44 liked a new flavor of Ben & Jerry's ice cream, we could never ask this question of all of the people in the U.S. who were in that age group. Such a research study would take way too much time to complete and the cost of doing that study would be prohibitive.

So, instead of testing *everyone* in the population, we take a sample of people in the population and use the results of this sample to estimate the mean of the entire population. This saves both time and money. When we use the results of a sample to estimate the population mean, this is called "*inferential statistics*" because we are inferring the population mean from the sample mean.

When we study a sample of people in business research, we know the size of our sample (n), the mean of our sample (\bar{X}), and the standard deviation of our sample (STDEV). We use these figures to estimate the population mean with a test called the “confidence interval about the mean.”

3.1.2 *Estimating the Lower Limit and the Upper Limit of the 95% Confidence Interval About the Mean*

The theoretical background of this test is beyond the scope of this book, and you can learn more about this test from studying any good statistics textbook (e.g. Levine 2011) but the basic ideas are as follows.

We assume that the population mean is somewhere in an interval which has a “lower limit” and an “upper limit” to it. We also assume in this book that we want to be “95% confident” that the population mean is inside this interval somewhere. So, we intend to make the following type of statement:

“We are 95% confident that the population mean in miles per gallon (mpg) for the Chevy Impala automobile is between 26.92 miles per gallon and 29.42 miles per gallon.”

If we want to create a billboard for this car that claims that this car gets 28 miles per gallon (mpg), we can do that because 28 is *inside the 95% confidence interval* in our research study in the above example. We do not know exactly what the population mean is, only that it is somewhere between 26.92 mpg and 29.42 mpg, and 28 is inside this interval.

But we are only 95% confident that the population mean is inside this interval, and 5% of the time we will be wrong in assuming that the population mean is 28 mpg.

But, for our purposes in business research, we are happy to be 95% confident that our assumption is accurate. We should also point out that 95% is an arbitrary level of confidence for our results. We could choose to be 80% confident, or 90% confident, or even 99% confident in our results if we wanted to do that. But, in this book, *we will always assume that we want to be 95% confident of our results*. That way, you will not have to guess on how confident you want to be in any of the problems in this book. We will always want to be 95% confident of our results in this book.

So how do we find the 95% confidence interval about the mean for our data?

In words, we will find this interval this way:

“Take the sample mean (\bar{X}), *and add to it* 1.96 times the standard error of the mean (s.e.) to get the upper limit of the confidence interval. Then, take the sample mean, *and subtract from it* 1.96 times the standard error of the mean to get the lower limit of the confidence interval.”

You will remember (See Sect. 1.3) that the standard error of the mean (s.e.) is found by dividing the standard deviation of our sample (STDEV) by the square root of our sample size, n .

In mathematical terms, the formula for the 95% confidence interval about the mean is:

$$\bar{X} \pm 1.96 \text{ s.e.} \quad (3.1)$$

Note that the “ \pm sign” stands for “plus or minus,” and this means that you first add 1.96 times the s.e. to the mean to get the upper limit of the confidence interval, and then subtract 1.96 times the s.e. from the mean to get the lower limit of the confidence interval. Also, the symbol 1.96 s.e. means that you multiply 1.96 times the standard error of the mean to get this part of the formula for the confidence interval.

Note: We will explain shortly where the number 1.96 came from.

Let’s try a simple example to illustrate this formula.

3.1.3 Estimating the Confidence Interval for the Chevy Impala in Miles Per Gallon

Let’s suppose that you asked owners of the Chevy Impala to keep track of their mileage and the number of gallons used for two tanks of gas. Let’s suppose that 49 owners did this, and that they average 27.83 miles per gallon (mpg) with a standard deviation of 3.01 mpg. The standard error (s.e.) would be 3.01 divided by the square root of 49 (i.e., 7) which gives a s.e. equal to 0.43.

The 95% confidence interval for these data would be:

$$27.83 \pm 1.96 (0.43)$$

The *upper limit of this confidence interval* uses the plus sign of the \pm sign in the formula. Therefore, the upper limit would be:

$$27.83 + 1.96 (0.43) = 27.83 + 0.84 = 28.67 \text{ mpg}$$

Similarly, the *lower limit of this confidence interval* uses the minus sign of the \pm sign in the formula. Therefore, the lower limit would be:

$$27.83 - 1.96 (0.43) = 27.83 - 0.84 = 26.99 \text{ mpg}$$

The result of our research study would, therefore, be the following:

“We are 95% confident that the population mean for the Chevy Impala is somewhere between 26.99 mpg and 28.67 mpg.”

If we were planning to create a billboard that claimed that this car got 28 mpg, we would be able to do that based on our data, since 28 is inside of this 95% confidence interval for the population mean.

You are probably asking yourself: “Where did that 1.96 in the formula come from?”

3.1.4 Where Did the Number “1.96” Come From?

A detailed mathematical answer to that question is beyond the scope of this book, but here is the basic idea.

We make an assumption that the data in the population are “normally distributed” in the sense that the population data would take the shape of a “normal curve” if we could test all of the people in the population. The normal curve looks like the outline of the Liberty Bell that sits in front of Independence Hall in Philadelphia, Pennsylvania. The normal curve is “symmetric” in the sense that if we cut it down the middle, and folded it over to one side, the half that we folded over would fit perfectly onto the half on the other side.

A discussion of integral calculus is beyond the scope of this book, but essentially we want to find the lower limit and the upper limit of the population data in the normal curve so that 95% of the area under this curve is between these two limits. *If we have more than 40 people in our research study*, the value of these limits is plus or minus 1.96 times the standard error of the mean (s.e.) of our sample. The number 1.96 times the s.e. of our sample gives us the upper limit and the lower limit of our confidence interval. If you want to learn more about this idea, you can consult a good statistics book (e.g. Salkind 2010).

The number 1.96 would change if we wanted to be confident of our results at a different level from 95% as long as we have more than 40 people in our research study.

For example:

1. If we wanted to be 80% confident of our results, this number would be 1.282.
2. If we wanted to be 90% confident of our results, this number would be 1.645.
3. If we wanted to be 99% confident of our results, this number would be 2.576.

But since we always want to be 95% confident of our results in this book, we will always use 1.96 in this book whenever we have more than 40 people in our research study.

By now, you are probably asking yourself: “Is this number in the confidence interval about the mean always 1.96?” The answer is: “No!”, and we will explain why this is true now.

3.1.5 Finding the Value for t in the Confidence Interval Formula

Objective: To find the value for t in the confidence interval formula

The correct formula for the confidence interval about the mean for different sample sizes is the following:

$$\bar{X} \pm t \text{ s.e.} \quad (3.2)$$

To use this formula, you find the sample mean, \bar{X} , and *add to it the value of t times the s.e. to get the upper limit* of this 95% confidence interval. Also, you take the sample mean, \bar{X} , and *subtract from it the value of t times the s.e. to get the lower limit* of this 95% confidence interval. And, you find the value of t in the table given in Appendix E of this book in the following way:

Objective: To find the value of t in the t-table in Appendix E

Before we get into an explanation of what is meant by “the value of t ,” let’s give you practice in finding the value of t by using the t-table in Appendix E.

Keep your finger on Appendix E as we explain how you need to “read” that table.

Since the test in this chapter is called the “confidence interval about the mean test,” you will use the first column on the left in Appendix E to find the critical value of t for your research study (note that this column is headed: “ sample size n ”).

To find the value of t , you go down this first column until you find the sample size in your research study, and then you go to the right and read the value of t for that sample size in the “critical t column” of the table (note that this column is the column that you would use for the 95% confidence interval about the mean).

For example, if you have 14 people in your research study, the value of t is 2.160.

If you have 26 people in your research study, the value of t is 2.060.

If you have more than 40 people in your research study, the value of t is always 1.96.

Note that the “critical t column” in Appendix E represents the value of t that you need to use to obtain to be 95% confident of your results as “significant” results.

Throughout this book, we are assuming that you want to be 95% confident in the results of your statistical tests. Therefore, the value for t in the t-table in Appendix E tells you which value you should use for t when you use the formula for the 95% confidence interval about the mean.

Now that you know how to find the value of t in the formula for the confidence interval about the mean, let’s explore how you find this confidence interval using Excel.

3.1.6 Using Excel's TINV Function to Find the Confidence Interval About the Mean

Objective: To use the TINV function in Excel to find the confidence interval about the mean

When you use Excel, the formulas for finding the confidence interval are:

$$\text{Lower limit: } = \bar{X} - TINV(1 - 0.95, n - 1) * s.e. \text{ (no spaces between these symbols)} \quad (3.3)$$

$$\text{Upper limit: } = \bar{X} + TINV(1 - 0.95, n - 1) * s.e. \text{ (no spaces between these symbols)} \quad (3.4)$$

Note that the “* symbol” in this formula tells Excel to use the multiplication step in the formula, and it stands for “times” in the way we talk about multiplication.

You will recall from Chap. 1 that n stands for the sample size, and so $n - 1$ stands for the sample size minus one.

You will also recall from Chap. 1 that the standard error of the mean, s.e., equals the STDEV divided by the square root of the sample size, n (See Sect. 1.3).

Let's try a sample problem using Excel to find the 95% confidence interval about the mean for a problem.

Suppose that General Motors wanted to claim that its Chevy Impala gets 28 miles per gallon (mpg), and that it wanted to advertise on a billboard in St. Louis at the Vandeventer entrance to Route 44: “The new Chevy Impala gets 28 miles to the gallon.” Let's call 28 mpg the “reference value” for this car.

Suppose that you work for Ford Motor Co. and that you want to check this claim to see if it holds up based on some research evidence. You decide to collect some data and to use a two-side 95% confidence interval about the mean to test your results:

3.1.7 Using Excel to Find the 95% Confidence Interval for a Car's mpg Claim

Objective: To analyze the data using a two-side 95% confidence interval about the mean

- C7: n
- C10: Mean
- C13: STDEV
- C16: s.e.
- C19: 95% confidence interval
- D21: Lower limit:
- D23: Upper limit: (see Fig. 3.2)

Chevy Impala				
Miles per gallon				
30.9				
24.5	n			
31.2				
28.7				
35.1	Mean			
29.0				
28.8				
23.1	STDEV			
31.0				
30.2				
28.4	s.e			
29.3				
24.2				
27.0	95% confidence interval			
26.7				
31.0			Lower limit:	
23.5				
29.4			Upper Limit:	
26.3				
27.5				
28.2				
28.4				
29.1				
21.9				
30.9				

Fig. 3.2 Example of Chevy Impala Format for the Confidence Interval About the Mean Labels

- B26: Draw a picture below this confidence interval
- B28: 26.92
- B29: lower (then right-align this word)
- B30: limit (then right-align this word)

- C28: ‘----- 28 -----28.17 ----- (note that you need to begin cell C28 with a *single quotation mark* (‘) to tell Excel that this is a *label*, and not a number)
- D28: ‘----- (notice the single quotation mark at the beginning)
- E28: ‘29.42 (note the single quotation mark)
- C29: ref. Mean
- C30: value
- E29: upper
- E30: limit
- B33: Conclusion:

Now, align the labels underneath the picture of the confidence interval so that they look like Fig. 3.3.

Chevy Impala				
Miles per gallon				
30.9				
24.5	n			
31.2				
28.7				
35.1	Mean			
29.0				
28.8				
23.1	STDEV			
31.0				
30.2				
28.4	s.e			
29.3				
24.2				
27.0	95% confidence interval			
26.7				
31.0			Lower limit:	
23.5				
29.4			Upper Limit:	
26.3				
27.5				
28.2	Draw a picture below this confidence interval			
28.4				
29.1	26.92 ----- 28 ----- 28.17 ----- 29.42			
21.9	lower	ref.	Mean	upper
30.9	limit	value		limit
	Conclusion:			

Fig. 3.3 Example of Drawing a Picture of a Confidence Interval About the Mean Result

Next, name the range of data from A6:A30 as: miles

D7: Use Excel to find the sample size

D10: Use Excel to find the mean

D13: Use Excel to find the STDEV

D16: Use Excel to find the s.e.

Now, you need to find the lower limit and the upper limit of the 95% confidence interval for this study.

We will use Excel's TINV function to do this. We will assume that you want to be 95% confident of your results.

F21: = D10 - TINV(1 - .95,24)*D16 (no spaces between)

Note that this TINV formula uses 24 since 24 is one less than the sample size of 25 (i.e., 24 is $n - 1$). Note that D10 is the mean, while D16 is the standard error of the mean. The above formula gives the *lower limit of the confidence interval*, 26.92.

F23: = D10 + TINV(1-.95,24)*D16 (no spaces between)

The above formula gives the *upper limit of the confidence interval*, 29.42.

Now, use number format (two decimal places) in your Excel spreadsheet for the mean, standard deviation, standard error of the mean, and for both the lower limit and the upper limit of your confidence interval. If you printed this spreadsheet now, the lower limit of the confidence interval (26.92) and the upper limit of the confidence interval (29.42) would “dribble over” onto a second printed page because the information on the spreadsheet is too large to fit onto one page in its present format.

So, you need to use Excel's “Scale to Fit” commands that we discussed in Chap. 2 (see Sect. 2.4) to reduce the size of the spreadsheet to 95% of its current size using the Page Layout/Scale to Fit function. Do that now, and notice that the dotted line to the right of 26.92 and 29.42 indicates that these numbers would now fit onto one page when the spreadsheet is printed out (see Fig. 3.4)

F21		fx =D10-TINV(1-0.95,24)*D16				
A	B	C	D	E	F	
31.2						
28.7						
35.1		Mean	28.17			
29.0						
28.8						
23.1		STDEV	3.03			
31.0						
30.2						
28.4		s.e	0.61			
29.3						
24.2						
27.0		95% confidence interval				
26.7						
31.0			Lower limit:	26.92		
23.5						
29.4			Upper Limit:	29.42		
26.3						
27.5						
28.2	Draw a picture below this confidence interval					
28.4						
29.1	26.92	28	28.17	29.42		
21.9	lower	ref.	Mean	upper		

Fig. 3.4 Result of Using the TINV Function to Find the Confidence Interval About the Mean

Note that you have drawn a picture of the 95% confidence interval beneath cell B26, including the lower limit, the upper limit, the mean, and the reference value of 28 mpg given in the claim that the company wants to make about the car’s miles per gallon performance.

Now, let’s write the conclusion to your research study on your spreadsheet:

- C33: Since the reference value of 28 is inside
- C34: the confidence interval, we accept that
- C35: the Chevy Impala does get 28 mpg.

Your research study accepted the claim that the Chevy Impala did get 28 miles per gallon. The average miles per gallon in your study was 28.17. (See Fig. 3.5)

Save your resulting spreadsheet as: **CHEVY7**

2. “If we increase our advertising budget by \$400,000 for our product, then our market share will go up by two points.”
3. “If we use this new method of teaching mathematics to ninth graders in algebra, then our math achievement scores will go up by 10%.”
4. “If we change the raw materials for this product, then our production cost per unit will decrease by 5%.”

A hypothesis, then, to a social science researcher is a “guess” about what we think is true in the real world. We can test these guesses using statistical formulas to see if our predictions come true in the real world.

So, in order to perform these statistical tests, we must first state our hypotheses so that we can test our results against our hypotheses to see if our hypotheses match reality.

So, how do we generate hypotheses in business?

3.2.1 Hypotheses Always Refer to the Population of People or Events That You Are Studying

The first step is to understand that our hypotheses always refer to the *population* of people under study.

For example, if we are interested in studying 18–24 year-olds in St. Louis as our target market, and we select a sample of people in this age group in St. Louis, depending on how we select our sample, we are hoping that our results of this study are useful in generalizing our findings to *all* 18–24 year-olds in St. Louis, and not just to the particular people in our sample.

The entire group of 18–24 year-olds in St. Louis would be the *population* that we are interested in studying, while the particular group of people in our study are called the *sample* from this population.

Since our sample sizes typically contain only a few people, we interested in the results of our sample *only insofar as the results of our sample can be “generalized” to the population in which we are really interested.*

That is why our hypotheses always refer to the population, and never to the sample of people in our study.

You will recall from Chap. 1 that we used the symbol: \bar{X} to refer to the mean of the sample we use in our research study (See Sect. 1.1).

We will use the symbol: μ (the Greek letter “mu”) to refer to the *population mean*.

In testing our hypotheses, we are trying to decide which one of two competing hypotheses *about the population mean* we should accept given our data set.

3.2.2 The Null Hypothesis and the Research (Alternative) Hypothesis

These two hypotheses are called the *null hypothesis* and the *research hypothesis*.

Statistics textbooks typically refer to the *null hypothesis* with the notation: H_0 .

The *research hypothesis* is typically referred to with the notation: H_1 , and it is sometimes called the *alternative hypothesis*.

Let’s explain first what is meant by the null hypothesis and the research hypothesis:

1. *The null hypothesis is what we accept as true unless we have compelling evidence that it is not true.*
2. *The research hypothesis is what we accept as true whenever we reject the null hypothesis as true.*

This is similar to our legal system in America where we assume that a supposed criminal is innocent until he or she is proven guilty in the eyes of a jury. Our null hypothesis is that this defendant is innocent, while the research hypothesis is that he or she is guilty.

In the great state of Missouri, every license plate has the state slogan: “Show me.” This means that people in Missouri think of themselves as not gullible enough to accept everything that someone says as true unless that person’s actions indicate the truth of his or her claim. In other words, people in Missouri believe strongly that a person’s actions speak much louder than that person’s words.

Since both the *null hypothesis* and the *research hypothesis* cannot both be true, the task of hypothesis testing using statistical formulas is to decide which one you will accept as true, and which one you will reject as true.

Sometimes in business research a series of rating scales is used to measure people’s attitudes toward a company, toward one of its products, or toward their intention-to-buy that company’s products. These rating scales are typically 5-point, 7-point, or 10-point scales, although other scale values are often used as well.

3.2.2.1 Determining the Null Hypothesis and the Research Hypothesis When Rating Scales Are Used

Here is a typical example of a 7-point scale in attitude research in customer satisfaction studies (see Fig. 3.6):



Fig. 3.6 Example of a Rating Scale Item for a New Car Purchase (Practical Example)

So, how do we decide what to use as the null hypothesis and the research hypothesis whenever rating scales are used?

Objective: To decide on the null hypothesis and the research hypothesis whenever rating scales are used.

In order to make this determination, we will use a simple rule:

Rule: Whenever rating scales are used, we will use the “middle” of the scale as the null hypothesis and the research hypothesis.

In the above example, since 4 is the number in the middle of the scale (i.e., three numbers are below it, and three numbers are above it), our hypotheses become:

Null hypothesis: $\mu = 4$

Research hypothesis: $\mu \neq 4$

In the above rating scale example, if the result of our statistical test for this one attitude scale item indicates that our population mean is “close to 4,” we say that we accept the null hypothesis that our new car purchase experience was neither positive nor negative.

In the above example, if the result of our statistical test indicates that the population mean is significantly different from 4, we reject the null hypothesis and accept the research hypothesis by stating either that:

“The new car purchase experience was significantly positive” (this is true whenever our sample mean is significantly greater than our expected population mean of 4).

or

“The new car purchase experience was significantly negative” (this is accepted as true whenever our sample mean is significantly less than our expected population mean of 4).

Both of these conclusions cannot be true. We accept one of the hypotheses as “true” based on the data set in our research study, and the other one as “not true” based on our data set.

The job of the business researcher, then, is to decide which of these two hypotheses, the null hypothesis or the research hypothesis, he or she will accept as true given the data set in the research study.

Let’s try some examples of rating scales so that you can practice figuring out what the null hypothesis and the research hypothesis are for each rating scale.

In the spaces in Fig. 3.7, write in the null hypothesis and the research hypothesis for the rating scales:

1. Webster University is an excellent university.									
	1	2	3	4	5				
	Strongly Disagree	Disagree	Undecided	Agree	Strongly Agree				
	Null hypothesis:			$\mu =$	_____				
	Research hypothesis:			$\mu \neq$	_____				
2. How would you rate the quality of teaching at Webster University?									
poor	1	2	3	4	5	6	7	excellent	
	Null hypothesis:			$\mu =$	_____				
	Research hypothesis:			$\mu \neq$	_____				
3. How would you rate the quality of the faculty at Webster University?									
1	2	3	4	5	6	7	8	9	10
very poor									very good
	Null hypothesis:			$\mu =$	_____				
	Research hypothesis:			$\mu \neq$	_____				

Fig. 3.7 Examples of Rating Scales for Determining the Null Hypothesis and the Research Hypothesis

How did you do?

Here are the answers to these three questions:

1. The null hypothesis is 3, and the research hypothesis is not equal to 3 on this 5-point scale (i.e. the “middle” of the scale is 3).
2. The null hypothesis is 4, and the research hypothesis is not equal to 4 on this 7-point scale (i.e., the “middle” of the scale is 4).
3. The null hypothesis is 5.5, and the research hypothesis is not equal to 5.5 on this 10-point scale (i.e., the “middle” of the scale is 5.5 since there are 5 numbers below 5.5 and 5 numbers above 5.5).

As another example, Holiday Inn Express in its Stay Smart Experience Survey uses 4-point scales where:

- 1 = Not So Good
- 2 = Average

3 = Very Good

4 = Great

On this scale, the null hypothesis is: $\mu = 2.5$ and the research hypothesis is: $\mu \neq 2.5$, because there are two numbers below 2.5, and two numbers above 2.5 on that rating scale.

Now, let's discuss the seven STEPS of hypothesis testing for using the confidence interval about the mean.

3.2.3 The Seven Steps for Hypothesis-Testing Using the Confidence Interval About the Mean

Objective: To learn the seven steps of hypothesis-testing using the confidence interval about the mean

There are seven basic steps of hypothesis-testing for this statistical test.

3.2.3.1 STEP 1: State the Null Hypothesis and the Research Hypothesis

If you are using numerical scales in your survey, you need to remember that these hypotheses refer to the “middle” of the numerical scale. For example, if you are using 7-point scales with 1 = poor and 7 = excellent, these hypotheses would refer to the middle of these scales and would be:

Null hypothesis H_0 : $\mu = 4$

Research hypothesis H_1 : $\mu \neq 4$

3.2.3.2 STEP 2: Select the Appropriate Statistical Test

In this chapter we are studying the confidence interval about the mean, and so we will select that test.

3.2.3.3 STEP 3: Calculate the Formula for the Statistical Test

You will recall (see Sect. 3.1.5) that the formula for the confidence interval about the mean is:

$$\bar{X} \pm t \text{ s.e.} \quad (3.2)$$

We discussed the procedure for computing this formula for the confidence interval about the mean using Excel earlier in this chapter, and the steps involved in using that formula are:

1. Use Excel's =COUNT function to find the sample size, n .
2. Use Excel's =AVERAGE function to find the sample mean, \bar{X} .
3. Use Excel's =STDEV function to find the standard deviation, STDEV.
4. Find the standard error of the mean (s.e.) by dividing the standard deviation (STDEV) by the square root of the sample size, n .
5. Use Excel's TINV function to find the lower limit of the confidence interval.
6. Use Excel's TINV function to find the upper limit of the confidence interval.

3.2.3.4 STEP 4: Draw a Picture of the Confidence Interval About the Mean, Including the Mean, the Lower Limit of the Interval, the Upper Limit of the Interval, and the Reference Value Given in the Null Hypothesis, H_0

(we will explain this step later in this chapter)

3.2.3.5 STEP 5: Decide on a Decision Rule

- (a) *If the reference value is inside the confidence interval, accept the null hypothesis, H_0*
- (b) *If the reference value is outside the confidence interval, reject the null hypothesis, H_0 , and accept the research hypothesis, H_1*

3.2.3.6 STEP 6: State the Result of Your Statistical Test

There are two possible results when you use the confidence interval about the mean, and only one of them can be accepted as "true." So your result would be one of the following:

Either: Since the reference value is inside the confidence interval, *we accept the null hypothesis, H_0*

Or: Since the reference value is outside the confidence interval, *we reject the null hypothesis, H_0 , and accept the research hypothesis, H_1*

3.2.3.7 STEP 7: State the Conclusion of Your Statistical Test in Plain English!

In practice, this is more difficult than it sounds because you are trying to summarize the result of your statistical test in simple English that is both concise and accurate so that someone who has never had a statistics course (such as your boss, perhaps)

can understand the conclusion of your test. This is a difficult task, and we will give you lots of practice doing this last and most important step throughout this book.

Objective: To write the conclusion of the confidence interval about the mean test

Let's set some basic rules for stating the conclusion of a hypothesis test.

Rule #1: *Whenever you reject H_0 and accept H_1 , you must use the word "significantly" in the conclusion to alert the reader that this test found an important result.*

Rule #2: *Create an outline in words of the "key terms" you want to include in your conclusion so that you do not forget to include some of them.*

Rule #3: *Write the conclusion in plain English so that the reader can understand it even if that reader has never taken a statistics course.*

Let's practice these rules using the Chevy Impala Excel spreadsheet that you created earlier in this chapter, but first we need to state the hypotheses for that car.

Since the billboard wants to claim that the Chevy Impala gets 28 miles per gallon, the hypotheses would be:

$$H_0 : \mu = 28 \text{ mpg}$$

$$H_1 : \mu \neq 28 \text{ mpg}$$

You will remember that the reference value of 28 mpg was inside the 95% confidence interval about the mean for your data, so we would accept H_0 for the Chevy Impala that the car does get 28 mpg.

Objective: To state the result when you accept H_0

Result: Since the reference value of 28 mpg is inside the confidence interval, we accept the null hypothesis, H_0

Let's try our three rules now:

Objective: To write the conclusion when you accept H_0

Rule #1: *Since the reference value was inside the confidence interval, we cannot use the word "significantly" in the conclusion. This is a basic rule we are using in this chapter for every problem.*

Rule #2: The key terms in the conclusion would be:

- Chevy Impala
- reference value of 28 mpg

Rule #3: The Chevy Impala did get 28 mpg.

The process of writing the conclusion when you accept H_0 is relatively straightforward since you put into words what you said when you wrote the null hypothesis.

However, the process of stating the conclusion when you reject H_0 and accept H_1 is more difficult, so let's practice writing that type of conclusion with three practice case examples:

Objective: To write the result and conclusion when you reject H_0

CASE #1: Suppose that an ad in *Business Week* claimed that the Ford Escape Hybrid got 34 miles per gallon. The hypotheses would be:

$$H_0 : \mu = 34 \text{ mpg}$$

$$H_1 : \mu \neq 34 \text{ mpg}$$

Suppose that your research yields the following confidence interval:

30	31	32	34
lower limit	Mean	upper limit	Ref. Value

Result: Since the reference value is outside the confidence interval, we reject the null hypothesis and accept the research hypothesis

The three rules for stating the conclusion would be:

Rule #1: We must include the word “significantly” since the reference value of 34 is outside the confidence interval.

Rule #2: The key terms would be:

- Ford Escape Hybrid
- significantly
- either “more than” or “less than”
- and probably closer to

Rule #3: The Ford Escape Hybrid got significantly less than 34 mpg, and it was probably closer to 31 mpg.

Note that this conclusion says that the mpg was less than 34 mpg because the sample mean was only 31 mpg. Note, also, that when you find a significant result by rejecting the null hypothesis, *it is not sufficient to say only: “significantly less than 34 mpg,”* because that does not tell the reader “how much less than 34 mpg” the sample mean was from 34 mpg. To make the conclusion clear, you need to add: “probably closer to 31 mpg” since the sample mean was only 31 mpg.

CASE #2: Suppose that you have been hired as a consultant by the St. Louis Symphony Orchestra (SLSO) to analyze the data from an Internet survey of attendees for a concert in Powell Symphony Hall in St. Louis last month. You have decided to practice your data analysis skills on Question #7 given in Fig. 3.8:

Question #7:	"Overall, how satisfied have you been with your experience(s) at SLSO concerts?"						
	1	2	3	4	5	6	7
	Extremely dissatisfied						Extremely satisfied

Fig. 3.8 Example of a Survey Item Used by the St. Louis Symphony Orchestra (SLSO)

The hypotheses for this one item would be:

$$H_0 : \mu = 4$$

$$H_1 : \mu \neq 4$$

Essentially, the null hypothesis equal to 4 states that if the obtained mean score for this question is not significantly different from 4 on the rating scale, then attendees, overall, were neither satisfied nor dissatisfied with their SLSO concerts.

Suppose that your analysis produced the following confidence interval for this item on the survey.

1.8	2.8	3.8	4
lower	Mean	upper	Ref.
limit		limit	Value

Result: Since the reference value is outside the confidence interval, we reject the null hypothesis and accept the research hypothesis.

Rule #1: You must include the word “significantly” since the reference value is outside the confidence interval

Rule #2: The key terms would be:

- attendees
- SLSO Internet survey
- significantly
- last month
- either satisfied or dissatisfied (since the result is significant)
- experiences at concerts
- overall

Rule #3: Attendees were significantly dissatisfied, overall, on last month’s Internet survey with their experiences at concerts of the SLSO.

Note that you need to use the word “dissatisfied” since the sample mean of 2.8 was on the dissatisfied side of the middle of the rating scale.

CASE #3: Suppose that Marriott Hotel at the St. Louis Airport location had the results of one item in its Guest Satisfaction Survey from last week's customers that was the following (see Fig. 3.9):

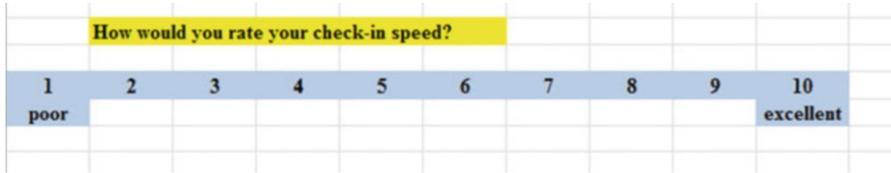


Fig. 3.9 Example of a Survey Item from Marriott Hotels

This item would have the following hypotheses:

$$H_0 : \mu = 5.5$$

$$H_1 : \mu \neq 5.5$$

Suppose that your research produced the following confidence interval for this item on the survey:

5.5	5.7	5.8	5.9
Ref.	lower	Mean	upper
Value	limit		limit

Result: Since the reference value is outside the confidence interval, we reject the null hypothesis and accept the research hypothesis

The three rules for stating the conclusion would be:

Rule #1: You must include the word “significantly” since the reference value is outside the confidence interval

Rule #2: The key terms would be:

- Marriott Hotel
- St. Louis Airport
- significantly
- check-in speed
- survey
- last week
- customers
- either “positive” or “negative” (we will explain this)

Rule #3: Customers at the St. Louis Airport Marriott Hotel last week rated their check-in speed in a survey as significantly positive.

Note two important things about this conclusion above: (1) people when speaking English do not normally say “significantly excellent” since something is either excellent or is not excellent without any modifier, and (2) since the mean rating of the check-in speed (5.8) was significantly greater than 5.5 on the positive side of the scale, we would say “significantly positive” to indicate this fact.

The three practice problems at the end of this chapter will give you additional practice in stating the conclusion of your result, and this book will include many more examples that will help you to write a clear and accurate conclusion to your research findings.

3.3 Alternative Ways to Summarize the Result of a Hypothesis Test

It is important for you to understand that in this book we are summarizing an hypothesis test in one of two ways: (1) We accept the null hypothesis, or (2) We reject the null hypothesis and accept the research hypothesis. We are consistent in the use of these words so that you can understand the concept underlying hypothesis testing.

However, there are many other ways to summarize the result of an hypothesis test, and all of them are correct theoretically, even though the terminology differs. If you are taking a course with a professor who wants you to summarize the results of a statistical test of hypotheses in language which is different from the language we are using in this book, do not panic! If you understand the concept of hypothesis testing as described in this book, you can then translate your understanding to use the terms that your professor wants you to use to reach the same conclusion to the hypothesis test.

Statisticians and professors of business statistics all have their own language that they like to use to summarize the results of an hypothesis test. There is no one set of words that these statisticians and professors will ever agree on, and so we have chosen the one that we believe to be easier to understand in terms of the concept of hypothesis testing.

To convince you that there are many ways to summarize the results of an hypothesis test, we present the following quotes from prominent statistics and research books to give you an idea of the different ways that are possible.

3.3.1 Different Ways to Accept the Null Hypothesis

The following quotes are typical of the language used in statistics and research books when the null hypothesis is accepted:

- “The null hypothesis is not rejected.” (Black 2010, p. 310)
- “The null hypothesis cannot be rejected.” (McDaniel and Gates 2010, p. 545)
- “The null hypothesis . . . claims that there is no difference between groups.” (Salkind 2010, p. 193)
- “The difference is not statistically significant.” (McDaniel and Gates 2010, p. 545)
- “... the obtained value is not extreme enough for us to say that the difference between Groups 1 and 2 occurred by anything other than chance.” (Salkind 2010, p. 225)
- “If we do not reject the null hypothesis, we conclude that there is not enough statistical evidence to infer that the alternative (hypothesis) is true.” (Keller 2009, p. 358)
- “The research hypothesis is not supported.” (Zikmund and Babin 2010, p. 552)

3.3.2 *Different Ways to Reject the Null Hypothesis*

The following quotes are typical of the quotes used in statistics and research books when the null hypothesis is rejected:

- “The null hypothesis is rejected.” (McDaniel and Gates 2010, p. 546)
- “If we reject the null hypothesis, we conclude that there is enough statistical evidence to infer that the alternative hypothesis is true.” (Keller 2009, p. 358)
- “If the test statistic’s value is inconsistent with the null hypothesis, we reject the null hypothesis and infer that the alternative hypothesis is true.” (Keller 2009, p. 348)
- “Because the observed value . . . is greater than the critical value . . . , the decision is to reject the null hypothesis.” (Black 2010, p. 359)
- “If the obtained value is more extreme than the critical value, the null hypothesis cannot be accepted.” (Salkind 2010, p. 243)
- “The critical t-value . . . must be surpassed by the observed t-value if the hypothesis test is to be statistically significant” (Zikmund and Babin 2010, p. 567)
- “The calculated test statistic . . . exceeds the upper boundary and falls into this rejection region. The null hypothesis is rejected.” (Weiers 2011, p. 330)

You should note that all of the above quotes are used by statisticians and professors when discussing the results of an hypothesis test, and so you should not be surprised if someone asks you to summarize the results of a statistical test using a different language than the one we are using in this book.

3.4 End-of-Chapter Practice Problems

1. Suppose that you have been asked by the manager of the *St. Louis Post-Dispatch* to analyze the data from a recent survey of past subscribers who have cancelled their newspaper subscription in the past 3 months. A random sample of this group was called by phone and asked a series of questions about the newspaper. The hypothetical data for survey question #4 appear in Fig. 3.10:

St. Louis Post-Dispatch Phone Survey				
Question #4: "How much would you be willing to pay per week for a six-month weekday/weekend subscription to the Post-Dispatch?"				
	Subscription Price (\$)			
	4.15			
	3.75			
	3.80			
	4.10			
	3.60			
	3.60			
	3.65			
	4.40			
	3.15			
	4.00			
	3.75			
	4.00			
	3.25			
	3.75			
	3.30			
	3.75			
	3.65			
	4.00			
	4.10			
	3.90			
	3.50			
	3.75			

Fig. 3.10 Worksheet Data for Chap. 3: Practice Problem #1

Suppose, further, that top management wants to charge \$3.80 for this new subscription price. Is this a reasonable price to charge based on the results of this survey question? (Hint: \$3.80 is the null hypothesis for this price.)

- (a) To the right of this table, use Excel to find the sample size, mean, standard deviation, and standard error of the mean for the price figures. Label your answers. Use currency format (two decimal places) for the mean, standard deviation, and standard error of the mean.
- (b) Enter the null hypothesis and the research hypothesis onto your spreadsheet.

- (c) Use Excel's TINV function to find the 95% confidence interval about the mean for these figures. Label your answers. Use currency format (two decimal places).
 - (d) Enter your *result* onto your spreadsheet.
 - (e) Enter your *conclusion in plain English* onto your spreadsheet.
 - (f) Print the final spreadsheet to fit onto one page (if you need help remembering how to do this, see the objectives at the end of Chap. 2 in Sect. 2.4)
 - (g) On your printout, draw a diagram of this 95% confidence interval by hand
 - (h) Save the file as: POST9
2. Suppose that you have been asked by the Human Resources department (HR) at your company to analyze the data from a recent "morale survey" of its managers to find out how managers think about working at your company. You want to test out your Excel skills on a small sample of managers with one item from the survey. You select a random sample of managers and the hypothetical data from Item #24 are given in Fig. 3.11.

HUMAN RESOURCES DEPARTMENT						
MORALE SURVEY OF MANAGERS						
Item #24: "How would you rate the quality of leadership shown by top management in this company?"						
1	2	3	4	5	6	7
very weak						very strong
			Rating			
			5			
			6			
			3			
			4			
			7			
			2			
			3			
			4			
			2			
			5			
			3			
			4			
			2			
			3			
			6			
			5			
			7			
			4			
			6			
			4			
			3			
			4			
			2			
			3			
			5			
			4			

Fig. 3.11 Worksheet Data for Chap. 3: Practice Problem #2

Create an Excel spreadsheet with these data.

- (a) Use Excel to the right of the table to find the sample size, mean, standard deviation, and standard error of the mean for these data. Label your answers, and use two decimal places for the mean, standard deviation, and standard error of the mean
 - (b) Enter the null hypothesis and the research hypothesis for this item on your spreadsheet.
 - (c) Use Excel's TINV function to find the 95% confidence interval about the mean for these data. Label your answers on your spreadsheet. Use two decimal places for the lower limit and the upper limit of the confidence interval.
 - (d) Enter the *result* of the test on your spreadsheet.
 - (e) Enter the *conclusion* of the test in plain English on your spreadsheet.
 - (f) Print your final spreadsheet so that it fits onto one page (if you need help remembering how to do this, see the objectives at the end of Chap. 2 in Sect. 2.4).
 - (g) Draw a picture of the confidence interval, including the reference value, onto your spreadsheet.
 - (h) Save the final spreadsheet as: top8
3. The American Advertising Federation (AAF) was established in 1905 and is the oldest national advertising trade association in the US. Its membership includes more than 40,000 advertising professionals who are part of the organization's corporate membership. AAF also sponsors an annual Student Advertising Career Conference which consists of 2 days of networking and seminar presentations of current trends in the advertising industry. Suppose that you have been asked to analyze the data from the survey that was emailed one week after the conference ended to students who attended this year's conference. The survey contained 15 items, and Item #15 is given in Fig. 3.12. You want to make sure that you can analyze the data correctly, so you have created some hypothetical data for this one item to test your Excel skills.

Student Advertising Career Conference							
Survey							
Item #15:	How likely are you to recommend to other advertising students that they attend next year's AAF Student Advertising Career Conference?						
	1	2	3	4	5	6	7
	Very Unlikely						Very Likely
		RATING					
		5					
		6					
		4					
		7					
		5					
		6					
		4					
		3					
		1					
		2					
		5					
		6					
		7					
		6					
		7					
		6					
		5					
		3					
		4					

Fig. 3.12 Worksheet Data for Chap. 3: Practice Problem #3

Create an Excel spreadsheet with these data.

- (a) To the right of this table, use Excel to find the sample size, mean, standard deviation, and standard error of the mean for the item ratings. Label your answers. Use number format (two decimal places) for the mean, standard deviation, and standard error of the mean.
- (b) Enter the null hypothesis and the research hypothesis onto your spreadsheet.

- (c) Use Excel's TINV function to find the 95% confidence interval about the mean for these figures. Label your answers. Use number format (two decimal places).
- (d) Enter your *result* onto your spreadsheet.
- (e) Enter your *conclusion in plain English* onto your spreadsheet.
- (f) Print the final spreadsheet to fit onto one page (if you need help remembering how to do this, see the objectives at the end of Chap. 2 in Sect. 2.4)
- (g) On your printout, draw a diagram of this 95% confidence interval by hand
- (h) Save the file as: AAF4

References

- Black, K. Business Statistics: for Contemporary Decision Making (6th ed.). Hoboken, NJ: John Wiley & Sons, Inc., 2010.
- Keller, G. Statistics for Management and Economics (8th ed.). Mason, OH: South-Western Cengage learning, 2009.
- Levine, D.M. Statistics for Managers using Microsoft Excel (6th ed.). Boston, MA: Prentice Hall/Pearson, 2011.
- McDaniel, C. and Gates, R. Marketing Research (8th ed.). Hoboken, NJ: John Wiley & Sons, Inc., 2010.
- Salkind, N.J. Statistics for People Who (think they) Hate Statistics (2nd Excel 2007 ed.). Los Angeles, CA: Sage Publications, 2010.
- Weiers, R.M. Introduction to Business Statistics (7th ed.). Mason, OH: South-Western Cengage Learning, 2011.
- Zikmund, W.G. and Babin, B.J. Exploring Marketing Research (10th ed.). Mason, OH: South-Western Cengage learning, 2010.

Chapter 4

One-Group t-Test for the Mean



In this chapter, you will learn how to use one of the most popular and most helpful statistical tests in business research: the one-group t-test for the mean.

The formula for the one-group t-test is as follows:

$$t = \frac{\bar{X} - \mu}{S_{\bar{X}}} \text{ where} \tag{4.1}$$

$$\text{s.e.} = S_{\bar{X}} = \frac{S}{\sqrt{n}} \tag{4.2}$$

This formula asks you to take the mean (\bar{X}) and subtract the population mean (μ) from it, and then divide the answer by the standard error of the mean (s.e.). The standard error of the mean equals the standard deviation divided by the square root of n (the sample size).

Let's discuss the seven STEPS of hypothesis testing using the one-group t-test so that you can understand how this test is used.

4.1 The Seven STEPS for Hypothesis-Testing Using the One-Group t-Test

Objective: To learn the seven steps of hypothesis-testing using the one-group t-test

Before you can try out your Excel skills on the one-group t-test, you need to learn the basic steps of hypothesis-testing for this statistical test. There are seven steps in this process:

4.1.1 STEP 1: State the Null Hypothesis and the Research Hypothesis

If you are using numerical scales in your survey, you need to remember that these hypotheses refer to the “middle” of the numerical scale. For example, if you are using 7-point scales with 1 = poor and 7 = excellent, these hypotheses would refer to the middle of these scales and would be:

Null hypothesis H_0 : $\mu = 4$

Research hypothesis H_1 : $\mu \neq 4$

As a second example, suppose that you worked for Honda Motor Company and that you wanted to place a magazine ad that claimed that the new Honda Fit got 35 miles per gallon (mpg). The hypotheses for testing this claim on actual data would be:

H_0 : $\mu = 35$ mpg

H_1 : $\mu \neq 35$ mpg

4.1.2 STEP 2: Select the Appropriate Statistical Test

In this chapter we will be studying the one-group t-test, and so we will select that test.

4.1.3 STEP 3: Decide on a Decision Rule for the One-Group t-Test

- (a) If the absolute value of t is less than the critical value of t , accept the null hypothesis.
- (b) If the absolute value of t is greater than the critical value of t , reject the null hypothesis and accept the research hypothesis.

You are probably saying to yourself: “That sounds fine, but how do I find the absolute value of t ?”

4.1.3.1 Finding the Absolute Value of a Number

To do that, we need another objective:

Objective: To find the absolute value of a number

If you took a basic algebra course in high school, you may remember the concept of “absolute value.” In mathematical terms, the absolute value of any number is *always* that number expressed as a positive number.

For example, the absolute value of 2.35 is +2.35.

And the absolute value of minus 2.35 (i.e. -2.35) is also +2.35.

This becomes important when you are using the t-table in Appendix E of this book. We will discuss this table later when we get to Step 5 of the one-group t-test where we explain how to find the critical value of t using Appendix E.

4.1.4 STEP 4: Calculate the Formula for the One-Group t-Test

Objective: To learn how to use the formula for the one-group t-test

The formula for the one-group t-test is as follows:

$$t = \frac{\bar{X} - \mu}{S_{\bar{X}}} \text{ where} \quad (4.1)$$

$$\text{s.e.} = S_{\bar{X}} = \frac{S}{\sqrt{n}} \quad (4.2)$$

This formula makes the following assumptions about the data (Foster et al. 1998): (1) The data are independent of each other (i.e., each person receives only one score), (2) the *population* of the data is normally distributed, and (3) the data have a constant variance (note that the standard deviation is the square root of the variance).

To use this formula, you need to follow these steps:

1. Take the sample mean in your research study and subtract the population mean μ from it (remember that the population mean for a study involving numerical rating scales is the “middle” number in the scale).
2. Then take your answer from the above step, and divide your answer by the standard error of the mean for your research study (you will remember that you learned how to find the standard error of the mean in Chap. 1; to find the standard error of the mean, just take the standard deviation of your research study and divide it by the square root of n , where n is the number of people used in your research study).
3. The number you get after you complete the above step is the value for t that results when you use the formula stated above.

4.1.5 *STEP 5: Find the Critical Value of t in the t-Table in Appendix E*

Objective: To find the critical value of t in the t-table in Appendix E

Before we get into an explanation of what is meant by “the critical value of t,” let’s give you practice in finding the critical value of t by using the t-table in Appendix E.

Keep your finger on Appendix E as we explain how you need to “read” that table.

Since the test in this chapter is called the “one-group t-test,” you will use the first column on the left in Appendix E to find the critical value of t for your research study (note that this column is headed: “sample size n”).

To find the critical value of t, you go down this first column until you find the sample size in your research study, and then you go to the right and read the critical value of t for that sample size in the critical t column in the table (note that *this is the column that you would use for both the one-group t-test and the 95% confidence interval about the mean*).

For example, if you have 27 people in your research study, the critical value of t is 2.056.

If you have 38 people in your research study, the critical value of t is 2.026.

If you have more than 40 people in your research study, the critical value of t is always 1.96.

Note that the “critical t column” in Appendix E represents the value of t that you need to obtain to be 95% confident of your results as “significant” results.

The critical value of t is the value that tells you whether or not you have found a “significant result” in your statistical test.

The t-table in Appendix E represents a series of “bell-shaped normal curves” (they are called bell-shaped because they look like the outline of the Liberty Bell that you can see in Philadelphia outside of Independence Hall).

The “middle” of these normal curves is treated as if it were zero point on the x-axis (the technical explanation of this fact is beyond the scope of this book, but any good statistics book (e.g. Zikmund and Babin 2010) will explain this concept to you if you are interested in learning more about it).

Thus, values of t that are to the right of this zero point are positive values that use a plus sign before them, and values of t that are to the left of this zero point are negative values that use a minus sign before them. Thus, some values of t are positive, and some are negative.

However, every statistics book that includes a t-table only reprints the *positive* side of the t-curves because the negative side is the mirror image of the positive side; this means that the negative side contains the exact same numbers as the positive side, but the negative numbers all have a minus sign in front of them.

Therefore, to use the t-table in Appendix E, you need to *take the absolute value of the t-value you found when you use the t-test formula* since the t-table in Appendix E only has the values of t that are the positive values for t.

Throughout this book, we are assuming that you want to be 95% confident in the results of your statistical tests. Therefore, the value for t in the t-table in Appendix E tells you whether or not the t-value you obtained when you used the formula for the one-group t-test is within the 95% interval of the t-curve range which that t-value would be expected to occur with 95% confidence.

If the t-value you obtained when you used the formula for the one-group t-test is *inside* of the 95% confidence range, we say that the result you found is *not significant* (note that this is equivalent to *accepting the null hypothesis!*).

If the t-value you found when you used the formula for the one-group t-test is *outside* of this 95% confidence range, we say that you have found a *significant result* that would be expected to occur less than 5% of the time (note that this is equivalent to *rejecting the null hypothesis and accepting the research hypothesis*).

4.1.6 STEP 6: State the Result of Your Statistical Test

There are two possible results when you use the one-group t-test, and only one of them can be accepted as “true.”

Either: Since the absolute value of t that you found in the t-test formula is *less than the critical value of t* in Appendix E, you accept the null hypothesis.

Or: Since the absolute value of t that you found in the t-test formula is *greater than the critical value of t* in Appendix E, you reject the null hypothesis, and accept the research hypothesis.

4.1.7 STEP 7: State the Conclusion of Your Statistical Test in Plain English!

In practice, this is more difficult than it sounds because you are trying to summarize the result of your statistical test in simple English that is both concise and accurate so that someone who has never had a statistics course (such as your boss, perhaps) can understand the result of your test. This is a difficult task, and we will give you lots of practice doing this last and most important step throughout this book.

If you have read this far, you are ready to sit down at your computer and perform the one-group t-test using Excel on some hypothetical data from the Guest Satisfaction Survey used by Marriott Hotels.

Let’s give this a try.

4.2 One-Group t-Test for the Mean

Suppose that you have been hired as a statistical consultant by Marriott Hotel in St. Louis to analyze the data from a Guest Satisfaction survey that they give to all customers to determine the degree of satisfaction of these customers for various activities of the hotel.

The survey contains a number of items, but suppose item #7 is the one in Fig. 4.1:

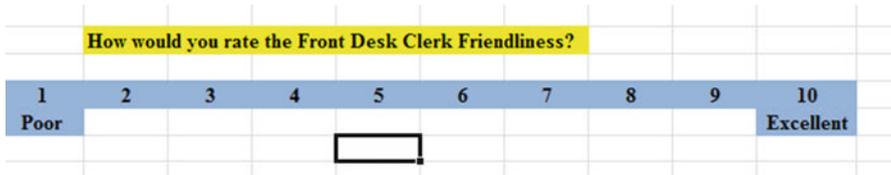


Fig. 4.1 Sample Survey Item for Marriot Hotel (Practical Example)

Suppose further, that you have decided to analyze the data from last week's customers using the one-group t-test.

Important note: You would need to use this test for each of the survey items separately.

Suppose that the hypothetical data for Item #7 from last week at the St. Louis Marriott Hotel were based on a sample size of 124 guests who had a mean score on this item of 6.58 and a standard deviation on this item of 2.44.

Objective: To analyze the data for each question separately using the one-group t-test for each survey item.

Create an Excel spreadsheet with the following information:

B11: Null hypothesis:

B14: Research hypothesis:

Note: Remember that when you are using a rating scale item, both the null hypothesis and the research hypothesis refer to the "middle of the scale." In the 10-point scale in this example, the middle of the scale is 5.5 since five numbers are below 5.5 (i.e., 1–5) and five numbers are above 5.5 (i.e. 6–10). Therefore, the hypotheses for this rating scale item are:

H₀ : $\mu = 5.5$

H₁ : $\mu \neq 5.5$

B17: n

B20: mean

B23: STDEV

B26: s.e.

- B29: critical t
- B32: t-test
- B36: Result:
- B41: Conclusion:

Now, use Excel:

- D17: enter the sample size
- D20: enter the mean
- D23: enter the STDEV (see Fig. 4.2)

Fig. 4.2 Basic Data Table for Front Desk Clerk Friendliness

Null hypothesis:				
Research hypothesis:				
n		124		
mean		6.58		
STDEV		2.44		
s.e.				
critical t				
t-test				
Result:				
Conclusion:				

D26: compute the standard error using the formula in Chap. 1

D29: find the critical t value of t in the t-table in Appendix E

Now, enter the following formula in cell D32 to find the t-test result:

$$=(D20-5.5)/D26 \quad (\text{no spaces between})$$

This formula takes the sample mean (D20) and subtracts the population hypothesized mean of 5.5 from the sample mean, and THEN divides the answer by the standard error of the mean (D26). Note that you need to enter $D20 - 5.5$ with an open-parenthesis *before* D20 and a closed-parenthesis *after* 5.5 so that the *answer of 1.08 is THEN divided by the standard error of 0.22* to get the t-test result of 4.93.

Now, use two decimal places for both the s.e. and the t-test result (see Fig. 4.3).

Fig. 4.3 t-test Formula
Result for Front Desk Clerk
Friendliness

Null hypothesis:		
Research hypothesis:		
n		124
mean		6.58
STDEV		2.44
s.e.		0.22
critical t		1.96
t-test		4.93
Result:		
Conclusion:		

Now, write the following sentence in D36–D39 to summarize the result of the t-test:

- D36: Since the absolute value of t of 4.93 is
D37: greater than the critical t of 1.96, we
D38: reject the null hypothesis and accept
D39: the research hypothesis.

Lastly, write the following sentence in D41–D43 to summarize the conclusion of the result for Item #7 of the Marriott Guest Satisfaction Survey:

- D41: St. Louis Marriott Hotel guests rated the
D42: Front Desk Clerks as significantly
D43: friendly last week.

Save your file as: MARRIOTT3

Print the final spreadsheet so that it fits onto one page as given in Fig. 4.4. Enter the null hypothesis and the research hypothesis by hand on your spreadsheet

Null hypothesis:	$\mu = 5.5$			
Research hypothesis:	$\mu \neq 5.5$			
n	124			
mean	6.58			
STDEV	2.44			
s.e.	0.22			
critical t	1.96			
t-test	4.93			
Result:	Since the absolute value of t of 4.93 is greater than the critical t of 1.96, we reject the null hypothesis and accept the research hypothesis.			
Conclusion:	St. Louis Marriott Hotel guests rated the Front Desk Clerks as significantly friendly last week.			

Fig. 4.4 Final Spreadsheet for Front Desk Clerk Friendliness

IMPORTANT NOTE: It is important for you to understand that “technically” the above conclusion in statistical terms should state:

“St. Louis Marriott Hotel Guests rated the Front Desk Clerks as friendly last week, and this result was probably not obtained by chance.”

However, throughout this book, we are using the term “significantly” in writing the conclusion of statistical tests to alert

the reader that the result of the statistical test was probably not a chance finding, but instead of writing all of those words each time, we use the word “significantly” as a shorthand to the longer explanation. This makes it much easier for the reader to understand the conclusion when it is written “in plain English,” instead of technical, statistical language.

4.3 Can You Use Either the 95% Confidence Interval About the Mean OR the One-Group t-Test When Testing Hypotheses?

You are probably asking yourself:

“It sounds like you could use *either* the 95% confidence interval about the mean *or* the one-group t-test to analyze the results of the types of problems described so far in this book? Is this a correct statement?”

The answer is a resounding: “Yes!”

Both the confidence interval about the mean and the one-group t-test are used often in business research on the types of problems described so far in this book. *Both of these tests produce the same result and the same conclusion from the data set!*

Both of these tests are explained in this book because some managers prefer the confidence interval about the mean test, others prefer the one-group t-test, and still others prefer to use both tests on the same data to make their results and conclusions clearer to the reader of their research reports. Since we do not know which of these tests your manager prefers, we have explained both of them so that you are competent in the use of both tests in the analysis of statistical data.

Now, let’s try your Excel skills on the one-group t-test on these three problems at the end of this chapter.

4.4 End-of-Chapter Practice Problems

1. Subaru of America rates the customer satisfaction of its dealers on a weekly basis on its Purchase Experience Survey, and demands that dealers achieve a 93% satisfaction score, or the dealers are required to take additional training to improve their customer satisfaction scores. Suppose that you have selected a random sample of rating forms submitted by new car purchasers (either online or through the mail) for the St. Louis Subaru dealer from a recent week and that you have prepared the hypothetical table in Fig. 4.5 for Question #1d:

SUBARU Customer Satisfaction Survey							
Question #Id:	"The salesperson was knowledgeable about the Subaru model line."						
	1	2	3	4	5	6	7
	Completely Disagree						Completely Agree
		Rating					
		5					
		7					
		6					
		4					
		3					
		5					
		6					
		7					
		2					
		3					
		5					
		7					
		4					
		7					
		7					
		5					
		6					
		6					
		4					
		3					
		5					
		5					

Fig. 4.5 Worksheet Data for Chap. 4: Practice Problem #1

- (a) Write the null hypothesis and the research hypothesis on your spreadsheet
 - (b) Use Excel to find the sample size, mean, standard deviation, and standard error of the mean to the right of the data set. Use number format (two decimal places) for the mean, standard deviation, and standard error of the mean.
 - (c) Enter the critical t from the t-table in Appendix E onto your spreadsheet, and label it.
 - (d) Use Excel to compute the t-value for these data (use two decimal places) and label it on your spreadsheet
 - (e) Type the result on your spreadsheet, and then type the conclusion in plain English on your spreadsheet
 - (f) Save the file as: `subaru4`
2. Boston University (BU) in Boston, MA US offers a graduate program for an M.S. degree in Advertising. One of the courses in this program focuses on Advertising Management. Suppose that you have been hired as a consultant by BU to analyze the student evaluation data for this course from the previous semester, and that you have created some hypothetical data for Question #12 that appear in Fig. 4.6.

BOSTON UNIVERSITY M.S. IN ADVERTISING PROGRAM							
Course: Advertising Management							
Item #12: "How would you rate the instructor's ability to explain advertising concepts clearly?"							
1	2	3	4	5	6	7	
Poor							Excellent
RATING							
5							
6							
4							
7							
6							
5							
7							
6							
7							
5							
6							
7							
6							
7							
5							
6							
7							
6							
7							

Fig. 4.6 Worksheet Data for Chap. 4: Practice Problem #2

- (a) Write the null hypothesis and the research hypothesis on your spreadsheet
 - (b) Use Excel to find the sample size, mean, standard deviation, and standard error of the mean to the right of the data set. Use number format (two decimal places) for the mean, standard deviation, and standard error of the mean.
 - (c) Enter the critical t from the t-table in Appendix E onto your spreadsheet, and label it.
 - (d) Use Excel to compute the t-value for these data (use two decimal places) and label it on your spreadsheet
 - (e) Type the result on your spreadsheet, and then type the conclusion in plain English on your spreadsheet
 - (f) Save the file as: COURSE3
3. Suppose that you have been hired as a marketing consultant by the Missouri Botanical Garden and have been asked to re-design the Comment Card survey that they have been asking visitors to The Garden to fill out after their visit. The Garden has been using a 5-point rating scale with 1 = poor and 5 = excellent. Suppose, further, that you have convinced The Garden staff to change to a 9-point

scale with 1 = poor and 9 = excellent so that the data will have a larger standard deviation. The hypothetical results of a recent week for Question #10 of your revised survey appear in Fig. 4.7.

MISSOURI BOTANICAL GARDEN								
VISITOR SURVEY								
Item #10: "How would you rate the helpfulness of The Garden staff?"								
1	2	3	4	5	6	7	8	9
poor								excellent
		Rating						
		8						
		6						
		5						
		7						
		9						
		5						
		6						
		4						
		8						
		7						
		6						
		8						
		6						
		7						
		9						
		7						
		6						
		3						
		8						
		7						
		6						

Fig. 4.7 Worksheet Data for Chap. 4: Practice problem #3

- (a) Write the null hypothesis and the research hypothesis on your spreadsheet
- (b) Use Excel to find the sample size, mean, standard deviation, and standard error of the mean to the right of the data set. Use number format (two decimal places) for the mean, standard deviation, and standard error of the mean.
- (c) Enter the critical t from the t-table in Appendix E onto your spreadsheet, and label it.
- (d) Use Excel to compute the t-value for these data (use two decimal places) and label it on your spreadsheet
- (e) Type the result on your spreadsheet, and then type the conclusion in plain English on your spreadsheet
- (f) Save the file as: Garden5

References

- Zikmund, W.G. and Babin, B.J. Exploring Marketing Research (10th ed.) Mason, OH: South-Western Cengage Learning, 2010.
- Foster, D.P., Stine, R.A., Waterman, R.P. Basic Business Statistics: A Casebook. New York, NY: Springer-Verlag, 1998.

Chapter 5

Two-Group t-Test of the Difference of the Means for Independent Groups



Up until now in this book, you have been dealing with the situation in which you have had only one group of people in your research study and only one measurement “number” on each of these people. We will now change gears and deal with the situation in which you are measuring two groups of people instead of only one group of people.

Whenever you have two completely different groups of people (i.e., no one person is in both groups, but every person is measured on only one variable to produce one “number” for each person), we say that the two groups are “independent of one another” This chapter deals with just that situation and that is why it is called the two-group t-test for independent groups.

The assumptions underlying the two-group t-test are the following (Zikmund and Babin 2010): (1) both groups are sampled from a normal population, and (2) the variances of the two populations are approximately equal. Note that the standard deviation is merely the square root of the variance. (There are different formulas to use when each person is measured twice to create two groups of data, and this situation is called “dependent,” but those formulas are beyond the scope of this book.) This book only deals with two groups that are independent of one another so that no person is in both groups of data.

When you are testing for the difference between the means for two groups, it is important to remember that there are two different formulas that you need to use depending on the sample sizes of the two groups:

1. Use Formula #1 in this chapter when both of the groups have more than 30 people in them, and
2. Use Formula #2 in this chapter when either one group, or both groups, have sample sizes less than 30 people in them.

We will illustrate both of these situations in this chapter.

But, first, we need to understand the steps involved in hypothesis-testing when two groups of people are involved before we dive into the formulas for this test.

5.1 The Nine STEPS for Hypothesis-Testing Using the Two-Group t-Test

Objective: To learn the nine steps of hypothesis-testing using two groups of people and the two-group t-test

You will see that these steps parallel the steps used in the previous chapter that dealt with the one-group t-test, but there are some important differences between the steps that you need to understand clearly before we dive into the formulas for the two-group t-test.

5.1.1 *STEP 1: Name One Group, Group 1, and the Other Group, Group 2*

The formulas used in this chapter will use the numbers 1 and 2 to distinguish between the two groups. If you define which group is Group 1 and which group is Group 2, you can use these numbers in your computations without having to write out the names of the groups.

For example, if you are testing teenage boys on their preference for the taste of Coke or Pepsi, you could call the groups: “Coke” and “Pepsi.” but this would require your writing out the words “Coke” or “Pepsi” whenever you wanted to refer to one of these groups. If you call the Coke group, Group 1, and the Pepsi group, Group 2, this makes it much easier to refer to the groups because it saves you writing time.

As a second example, you could be comparing the test market results for Kansas City vs. Indianapolis, but if you had to write out the names of those cities whenever you wanted to refer to them, it would take you more time than it would if, instead, you named one city, Group 1, and the other city, Group 2.

Note, also, that it is completely arbitrary which group you call Group 1, and which Group you call Group 2. You will achieve the same result and the same conclusion from the formulas however you decide to define these two groups.

5.1.2 *STEP 2: Create a Table That Summarizes the Sample Size, Mean Score, and Standard Deviation of Each Group*

This step makes it easier for you to make sure that you are using the correct numbers in the formulas for the two-group t-test. If you get the numbers “mixed-up,” your entire formula work will be incorrect and you will botch the problem terribly.

For example, suppose that you tested teenage boys on their preference for the taste of Coke vs. Pepsi in which the boys were randomly assigned to taste just one of these brands and then rate its taste on a 100-point scale from 0 = poor to 100 = excellent. After the research study was completed, suppose that the Coke group had 52 boys in it, their mean taste rating was 55 with a standard deviation of 7, while the Pepsi group had 57 boys in it and their average taste rating was 64 with a standard deviation of 13.

The formulas for analyzing these data to determine if there was a significant difference in the taste rating for teenage boys for these two brands require you to use six numbers correctly in the formulas: the sample size, the mean, and the standard deviation of each of the two groups. All six of these numbers must be used correctly in the formulas if you are to analyze the data correctly.

If you create a table to summarize these data, a good example of the table, using both Step 1 and Step 2, would be the data presented in Fig. 5.1:

Fig. 5.1 Basic Table Format for the Two-group t-test

Group	n	Mean	STDEV
1 (name it)			
2 (name it)			

For example, if you decide to call Group 1 the Coke group and Group 2 the Pepsi group, the following table would place the six numbers from your research study into the proper calls of the table as in Fig. 5.2:

Fig. 5.2 Results of Entering the Data Needed for the Two-group t-test

Group	n	Mean	STDEV
1 (name it)	52	55	7
2 (name it)	57	64	13

You can now use the formulas for the two-group t-test with more confidence that the six numbers will be placed in the proper place in the formulas.

Note that you could just as easily call Group 1 the Pepsi group and Group 2 the Coke group; it makes no difference how you decide to name the two groups; this decision is up to you.

5.1.3 STEP 3: State the Null Hypothesis and the Research Hypothesis for the Two-Group t-Test

If you have completed Step 1 above, this step is very easy because the null hypothesis and the research hypothesis will always be stated in the same way for the two-group t-test. The null hypothesis states that the population means of the two groups are equal, while the research hypothesis states that the population means of the two groups are not equal. In notation format, this becomes:

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2$$

You can now see that this notation is much simpler than having to write out the names of the two groups in all of your formulas.

5.1.4 STEP 4: Select the Appropriate Statistical Test

Since this chapter deals with the situation in which you have two groups of people but only one measurement on each person in each group, we will use the two-group t-test throughout this chapter.

5.1.5 STEP 5: Decide on a Decision Rule for the Two-Group t-Test

The decision rule is exactly what it was in the previous chapter (see Sect. 4.1.3) when we dealt with the one-group t-test.

- (a) If the absolute value of t is less than the critical value of t , accept the null hypothesis.
- (b) If the absolute value of t is greater than the critical value of t , reject the null hypothesis and accept the research hypothesis.

Since you learned how to find the absolute value of t in the previous chapter (see Sect. 4.1.3.1), you can use that knowledge in this chapter.

5.1.6 STEP 6: Calculate the Formula for the Two-Group t-Test

Since we are using two different formulas in this chapter for the two-group t-test depending on the sample size of the people in the two groups, we will explain how to use those formulas later in this chapter.

5.1.7 STEP 7: Find the Critical Value of t in the t -Table in Appendix E

In the previous chapter where we were dealing with the one-group t -test, you found the critical value of t in the t -table in Appendix E by finding the sample size for the one group of people in the first column of the table, and then reading the critical value of t across from it on the right in the “critical t column” in the table (see Sect. 4.1.5). This process was fairly simple once you have had some practice in doing this step.

However, for the two-group t -test, the procedure for finding the critical value of t is more complicated because you have two different groups of people in your study, and they often have different sample sizes in each group.

To use Appendix E correctly in this chapter, you need to learn how to find the “degrees of freedom” for your study. We will discuss that process now.

5.1.7.1 Finding the Degrees of Freedom (df) for the Two-Group t -Test

Objective: To find the degrees of freedom for the two-group t -test and to use it to find the critical value of t in the t -table in Appendix E

The mathematical explanation of the concept of the “degrees of freedom” is beyond the scope of this book, but you can find out more about this concept by reading any good statistics book (e.g. Keller 2009). For our purposes, you can easily understand how to find the degrees of freedom and to use it to find the critical value of t in Appendix E. The formula for the degrees of freedom (df) is:

$$\text{degrees of freedom} = df = n_1 + n_2 - 2 \quad (5.1)$$

In other words, you add the sample size for Group 1 to the sample size for Group 2 and then subtract 2 from this total to get the number of degrees of freedom to use in Appendix E.

Take a look at Appendix E.

Instead of using the first column as we did in the one-group t -test that is based on the sample size, n , of one group of people, *we need to use the second-column of this table (df) to find the critical value of t for the two-group t -test.*

For example, if you had 13 people in Group 1 and 17 people in Group 2, the degrees of freedom would be: $13 + 17 - 2 = 28$, and the critical value of t would be 2.048 *since you look down the second column which contains the degrees of freedom until you come to the number 28, and then read 2.048 in the “critical t column” in the table to find the critical value of t when $df = 28$.*

As a second example, if you had 52 people in Group 1 and 57 people in Group 2, the degrees of freedom would be: $52 + 57 - 2 = 107$. When you go down the

second column in Appendix E for the degrees of freedom, you find that *once you go beyond the degrees of freedom equal to 39, the critical value of t is always 1.96*, and that is the value you would use for the critical t with this example.

5.1.8 STEP 8: State the Result of Your Statistical Test

The result follows the exact same result format that you found for the one-group t -test in the previous chapter (see Sect. 4.1.6):

Either: Since the absolute value of t that you found in the t -test formula is *less than the critical value of t* in Appendix E, you accept the null hypothesis.

Or: Since the absolute value of t that you found in the t -test formula is *greater than the critical value of t* in Appendix E, you reject the null hypothesis and accept the research hypothesis.

5.1.9 STEP 9: State the Conclusion of Your Statistical Test in Plain English!

Writing the conclusion for the two-group t -test is more difficult than writing the conclusion for the one-group t -test because you have to decide what the difference was between the two groups.

When you accept the null hypothesis, the conclusion is simple to write: “There is no difference between the two groups in the variable that was measured.”

But when you reject the null hypothesis and accept the research hypothesis, you need to be careful about writing the conclusion so that it is both accurate and concise.

Let’s give you some practice in writing the conclusion of a two-group t -test.

5.1.9.1 Writing the Conclusion of the Two-Group t -Test When You Accept the Null Hypothesis

Objective: To write the conclusion of the two-group t -test when you have accepted the null hypothesis.

Suppose that you have been hired as a statistical consultant by Marriott Hotel in St. Louis to analyze the data from a Guest Satisfaction Survey that they give to all customers to determine the degree of satisfaction of these customers for various activities of the hotel.

The survey contains a number of items, but suppose Item #7 is the one in Fig. 5.3:

How would you rate the Front Desk Clerk Friendliness?									
1	2	3	4	5	6	7	8	9	10
Poor									Excellent

Fig. 5.3 Marriott Hotel Guest Satisfaction Survey Item #7

Suppose further, that you have decided to analyze the data from last week’s customers comparing men and women using the two-group t-test.

Important note: You would need to use this test for each of the survey items separately.

Suppose that the hypothetical data for Item #7 from last week at the St. Louis Marriott Hotel were based on a sample size of 124 men who had a mean score on this item of 6.58 and a standard deviation on this item of 2.44. Suppose that you also had data from 86 women from last week who had a mean score of 6.45 with a standard deviation of 1.86.

We will explain later in this chapter how to produce the results of the two-group t-test using its formulas, but, for now, let’s “cut to the chase” and tell you that these data would produce the following in Fig. 5.4:

Fig. 5.4 Worksheet Data for Males vs. Females for the St. Louis Marriott Hotel for Accepting the Null Hypothesis

Group	n	Mean	STDEV
1 Males	124	6.58	2.44
2 Females	86	6.45	1.86

- degrees of freedom: 208
- critical t: 1.96 (in Appendix E)
- t-test formula: 0.44 (when you use your calculator!)
- Result: Since the absolute value of 0.44 is less than the critical t of 1.96, we accept the null hypothesis.
- Conclusion: There was no difference between male and female guests last week in their rating of the friendliness of the front-desk clerk at the St. Louis Marriott Hotel.

Now, let’s see what happens when you reject the null hypothesis (H_0) and accept the research hypothesis (H_1).

5.1.9.2 Writing the Conclusion of the Two-Group t-Test When You Reject the Null Hypothesis and Accept the Research Hypothesis

Objective: To write the conclusion of the two-group t-test when you have rejected the null hypothesis and accepted the research hypothesis

Let’s continue with this same example of the Marriott Hotel, but with the result that we reject the null hypothesis and accept the research hypothesis.

Let’s assume that this time you have data on 85 males from last week and their mean score on this question was 7.26 with a standard deviation of 2.35. Let’s further suppose that you also have data on 48 females from last week and their mean score on this question was 4.37 with a standard deviation of 3.26. Let Males = Group 1 and Females = Group 2.

Without going into the details of the formulas for the two-group t-test, these data would produce the following result and conclusion based on Fig. 5.5:

Fig. 5.5 Worksheet Data for St. Louis Marriott Hotel for Obtaining a Significant Difference between Males and Females

Group	n	Mean	STDEV
1 Males	85	7.26	2.35
2 Females	48	4.37	3.26

Null Hypothesis: $\mu_1 = \mu_2$
 Research Hypothesis: $\mu_1 \neq \mu_2$
 degrees of freedom: 131
 critical t: 1.96 (in Appendix E)
 t-test formula: 5.40 (when you use your calculator!)
 Result: Since the absolute value of 5.40 is greater than the critical t of 1.96, we reject the null hypothesis and accept the research hypothesis.

Now, you need to compare the ratings of the men and women to find out which group had the more positive rating of the friendliness of the front-desk clerk using the following rule:

Rule: To summarize the conclusion of the two-group t-test, just compare the means of the two groups, and be sure to use the word “significantly” in your conclusion if you rejected the null hypothesis and accepted the research hypothesis.

A good way to prepare to write the conclusion of the two-group t-test when you are using a rating scale is to place the mean scores of the two groups on a drawing of the scale so that you can visualize the difference of the mean scores. For example, for our Marriott Hotel example above, you would draw this “picture” of the scale in Fig. 5.6:

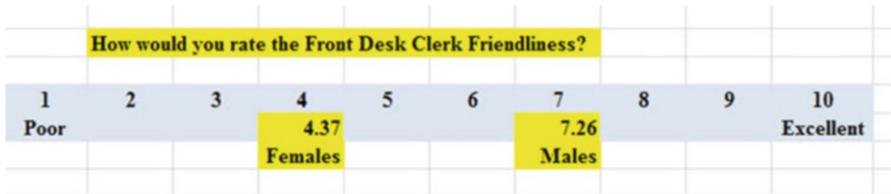


Fig. 5.6 Example of Drawing a “Picture” of the Means of the Two Groups on the Rating Scale

This drawing tells you visually that males had a higher positive rating than females on this item (7.26 vs. 4.37). *And, since you rejected the null hypothesis and accepted the research hypothesis, you know that you have found a significant difference between the two mean scores.*

So, our conclusion needs to contain the following key words:

- Male guests
- Female guests
- Marriott Hotel
- St. Louis
- last week
- significantly
- Front Desk Clerks
- more friendly *or* less friendly
- *either* (7.26 vs. 4.37) *or* (4.37 vs. 7.26)

We can use these key words to write the either of two conclusions which are *logically identical*:

Either: Male guests at the Marriott Hotel in St. Louis last week rated the Front Desk Clerks as significantly more friendly than Female guests (7.26 vs. 4.37).

Or: Female guests at the Marriott Hotel in St. Louis last week rated the Front Desk Clerks as significantly less friendly than Male guests (4.37 vs. 7.26).

Both of these conclusions are accurate, so you can decide which one you want to write. It is your choice.

Also, note that the mean scores in parentheses at the end of these conclusions must match the sequence of the two groups in your conclusion. For example, if you say that: “Male guests rated the Front Desk Clerks as significantly more friendly than Female guests,” the end of this conclusion should be: (7.26 vs. 4.37) since you mentioned Males first and Females second.

Alternately, if you wrote that: “Female guests rated the Front Desk Clerks as significantly less friendly than Male guests,” the end of this conclusion should be: (4.37 vs. 7.26) since you mentioned Females first and Males second.

Putting the two mean scores at the end of your conclusion saves the reader from having to turn back to the table in your research report to find these mean scores to see how far apart the mean scores were.

Now, let’s discuss FORMULA #1 that deals with the situation in which both groups have more than 30 people in them.

Objective: To use FORMULA #1 for the two-group t-test when both groups have a sample size greater than 30 people

5.2 Formula #1: Both Groups Have More Than 30 People in Them

The first formula we will discuss will be used when you have two groups of people with more than 30 people in each group and one measurement on each person in each group. This formula for the two-group t-test is:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}} \quad (5.2)$$

$$\text{where } S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \quad (5.3)$$

$$\text{and where degrees of freedom} = df = n_1 + n_2 - 2 \quad (5.1)$$

This formula looks daunting when you first see it, but let’s explain some of the parts of this formula:

We have explained the concept of “degrees of freedom” earlier in this chapter, and so you should be able to find the degrees of freedom needed for this formula in order to find the critical value of t in Appendix E.

In the previous chapter, *the formula for the one-group t-test was the following:*

$$t = \frac{\bar{X} - \mu}{S_{\bar{X}}} \quad (4.1)$$

$$\text{where s.e.} = S_{\bar{X}} = \frac{S}{\sqrt{n}} \quad (4.2)$$

For the one-group t-test, you found the mean score and subtracted the population mean from it, and then divided the result by the standard error of the mean (s.e.) to

get the result of the t-test. You then compared the t-test result to the critical value of t to see if you either accepted the null hypothesis, or rejected the null hypothesis and accepted the research hypothesis.

The two-group t-test requires a different formula because you have two groups of people, each with a mean score on some variable. You are trying to determine whether to accept the null hypothesis that the *population means of the two groups are equal* (in other words, there is no difference statistically between these two means), or whether the difference between the means of the two groups is “sufficiently large” that you would accept *that there is a significant difference* in the mean scores of the two groups.

The numerator of the two-group t-test asks you to find the difference of the means of the two groups:

$$\bar{X}_1 - \bar{X}_2 \quad (5.4)$$

The next step in the formula for the two-group t-test is to divide the answer you get when you subtract the two means by the standard error of the difference of the two means, and *this is a different standard error of the mean that you found for the one-group t-test because there are two means in the two-group t-test.*

The standard error of the mean when you have two groups of people is called the “standard error of the difference of the means” between the two groups. This formula looks less scary when you break it down into four steps:

1. Square the standard deviation of Group 1, and divide this result by the sample size for Group 1 (n_1).
2. Square the standard deviation of Group 2, and divide this result by the sample size for Group 2 (n_2).
3. Add the results of the above two steps to get a total score.
4. *Take the square root of this total score* to find the standard error of the difference of the means between the two groups, $S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$

This last step is the one that gives students the most difficulty when they are finding this standard error using their calculator, because they are in such a hurry to get to the answer that they forget to carry the square root sign down to the last step, and thus get a larger number than they should for the standard error.

5.2.1 An Example of Formula #1 for the Two-Group t-Test

Now, let’s use Formula #1 in a situation in which both groups have a sample size greater than 30 people.

Suppose that you have been hired by PepsiCo to do a taste test with teenage boys (ages 13–18) to determine if they like the taste of Pepsi the same as the taste of Coke. The boys are not told the brand name of the soft drink that they taste.

You select a group of boys in this age range, and randomly assign them to one of two groups: (1) Group 1 tastes Coke, and (2) Group 2 tastes Pepsi. Each group rates the taste of their soft drink on a 100-point scale using the following scale in Fig. 5.7:



Fig. 5.7 Example of a Rating Scale for a Soft Drink Taste Test (Practical Example)

Suppose you collect these ratings and determine (using your new Excel skills) that the 52 boys in the Coke group had a mean rating of 55 with a standard deviation of 7, while the 57 boys in the Pepsi group had a mean rating of 64 with a standard deviation of 13.

Note that the two-group t-test does not require that both groups have the same sample size. This is another way of saying that the two-group t-test is “robust” (a fancy term that statisticians like to use).

Your data then produce the following table in Fig. 5.8:

Fig. 5.8 Worksheet Data for Soft Drink Taste Test

Group	n	Mean	STDEV
1 Coke	52	55	7
2 Pepsi	57	64	13

Create an Excel spreadsheet, and enter the following information:

- B3: Group
- B4: 1 Coke
- B5: 2 Pepsi
- C3: n
- D3: Mean
- E3: STDEV
- C4: 52
- D4: 55
- E4: 7
- C5: 57
- D5: 64
- E5: 13

Now, widen column B so that it is twice as wide as column A, and center the six numbers and their labels in your table (see Fig. 5.9)

	A	B	C	D	E	F
1						
2						
3		Group	n	Mean	STDEV	
4		1 Coke	52	55	7	
5		2 Pepsi	57	64	13	
6						

Fig. 5.9 Results of Widening Column B and Centering the Numbers in the Cells

B8: Null hypothesis:

B10: Research hypothesis:

Since both groups have a sample size greater than 30, you need to use Formula #1 for the t-test for the difference of the means of the two groups.

Let’s “break this formula down into pieces” to reduce the chance of making a mistake.

B13: $STDEV1^2/n1$ (note that you square the standard deviation of Group 1, and then divide the result by the sample size of Group 1)

B16: $STDEV2^2/n2$

B19: $D13 + D16$

B22: s.e.

B25: critical t

B28: t-test

B31: Result:

B36: Conclusion: (see Fig. 5.10)

Fig. 5.10 Formula Labels for the Two-group t-test

Group	n	Mean	STDEV
1 Coke	52	55	7
2 Pepsi	57	64	13
Null hypothesis:			
Research hypothesis:			
STDEV1 squared / n1			
STDEV2 squared / n2			
D13 + D16			
s.e.			
critical t			
t-test			
Result:			
Conclusion:			

You now need to compute the values of the above formulas in the following cells:

- D13: the result of the formula needed to compute cell B13 (use two decimals)
- D16: the result of the formula needed to compute cell B16 (use two decimals)
- D19: the result of the formula needed to compute cell B19 (use two decimals)
- D22: =SQRT(D19) (use two decimals)

This formula should give you a standard error (s.e.) of 1.98.

D25: 1.96

(Since $df = n1 + n2 - 2$, this gives $df = 109 - 2 = 107$, and the critical t is, therefore, 1.96 in Appendix E.)

D28: $= (D4 - D5) / D22$ (use two decimals) (no spaces between)

This formula should give you a value for the t-test of: -4.55.

Next, check to see if you have rounded off all figures in D13: D28 to two decimal places (see Fig. 5.11).

Fig. 5.11 Results of the t-test Formula for the Soft Drink Taste Test

	A	B	C	D	E
11					
12					
13		STDEV1 squared / n1		0.94	
14					
15					
16		STDEV2 squared / n2		2.96	
17					
18					
19		D13 + D16		3.91	
20					
21					
22		s.e.		1.98	
23					
24					
25		critical t		1.96	
26					
27					
28		t-test		-4.55	
29					

Now, write the following sentence in D31 to D34 to summarize the result of the study:

- D31: Since the absolute value of -4.55
- D32: is greater than the critical t of
- D33: 1.96, we reject the null hypothesis
- D34: and accept the research hypothesis.

Finally, write the following sentence in D36 to D38 to summarize the conclusion of the study in plain English:

- D36: Teenage boys rated the taste of
- D37: Pepsi as significantly better than
- D38: the taste of Coke (64 vs. 55).

Save your file as: COKE4

Print this file so that it fits onto one page, and write by hand the null hypothesis and the research hypothesis on your printout.

The final spreadsheet appears in Fig. 5.12.

Group	n	Mean	STDEV
1 Coke	52	55	7
2 Pepsi	57	64	13
Null hypothesis:		$\mu_1 = \mu_2$	
Research hypothesis:		$\mu_1 \neq \mu_2$	
STDEV1 squared / n1		0.94	
STDEV2 squared / n2		2.96	
D13 + D16		3.91	
s.e.		1.98	
critical t		1.96	
t-test		-4.55	
Result:		Since the absolute value of - 4.55 is greater than the critical t of 1.96, we reject the null hypothesis and accept the research hypothesis.	
Conclusion:		Teenage boys rated the taste of Pepsi as significantly better than the taste of Coke (64 vs. 55)	

Fig. 5.12 Final Worksheet for the Coke vs. Pepsi Taste Test

Now, let's use the second formula for the two-group t-test which we use whenever either one group, or both groups, have less than 30 people in them.

Objective: To use Formula #2 for the two-group t-test when one or both groups have less than 30 people in them

Now, let’s look at the case when one or both groups have a sample size less than 30 people in them.

5.3 Formula #2: One or Both Groups Have Less Than 30 People in Them

Suppose that you work for the manufacturer of MP3 players and that you have been asked to do a pricing experiment to see if more units can be sold at a reduction in price.

Suppose, further, that you have randomly selected 7 wholesalers to purchase the product at the regular price, and they purchased a mean of 117.7 units with a standard deviation of 19.9 units.

In addition, you randomly selected a different group of 8 wholesalers to purchase the product at a 10% price cut, and they purchased a mean of 125.1 units with a standard deviation of 15.1 units. Let Regular Price = Group 1, and Reduced Price = Group 2.

You want to test to see if the two different prices produced a significant difference in the number of MP3 units sold.

You have decided to use the two-group t-test for independent samples, and the following data resulted in Fig. 5.13:

Null hypothesis: $\mu_1 = \mu_2$

Research hypothesis: $\mu_1 \neq \mu_2$

Group	n	Mean	STDEV
1 Regular Price	7	117.7	19.9
2 Reduced price	8	125.1	15.1

Fig. 5.13 Worksheet Data for Wholesaler Price Comparison (Practical Example)

Note: Since both groups have a sample size less than 30 people, you need to use Formula #2 in the following steps:

Create an Excel spreadsheet, and enter the following information:

- B3: Group
- B4: 1 Regular Price
- B5: 2 Reduced Price
- C3: n
- D3: Mean
- E3: STDEV

Now, widen column B so that it is three times as wide as column A.

To do this, click on B at the top left of your spreadsheet to highlight all of the cells in column B. Then, move the mouse pointer to the right end of the B cell until you get a “cross” sign; then, click on this cross sign and drag the sign to the right until you can read all of the words on your screen. Then, stop clicking!

- C4: 7
- D4: 117.7
- E4: 19.9
- C5: 8
- D5: 125.1
- E5: 15.1

Next, center the information in cells C3 to E5 by highlighting these cells and then using this step:

Click on the bottom line, second from the left icon, under “Alignment” at the top-center of Home

- B8: Null hypothesis:
- B10: Research hypothesis: (See Fig. 5.14)

	A	B	C	D	E	F	G
1							
2							
3		Group	n	Mean	STDEV		
4		1 Regular Price	7	117.7	19.9		
5		2 Reduced Price	8	125.1	15.1		
6							
7							
8		Null hypothesis:					
9							
10		Research hypothesis:					
11							

Fig. 5.14 Wholesaler Price Comparison Worksheet Data for Hypothesis Testing

Since both groups have a sample size less than 30, you need to use Formula #2 for the t-test for the difference of the means of two independent samples.

Formula #2 for the two-group t-test is the following:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}} \tag{5.2}$$

where $S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$ (5.5)

and where degrees of freedom = $df = n_1 + n_2 - 2$ (5.6)

This formula is complicated, and so it will reduce your chance of making a mistake in writing it if you “break it down into pieces” instead of trying to write the formula as one cell entry.

Now, enter these words on your spreadsheet:

- B13: $(n_1 - 1) \times \text{STDEV1 squared}$
- B16: $(n_2 - 1) \times \text{STDEV2 squared}$
- B19: $n_1 + n_2 - 2$
- B22: $1/n_1 + 1/n_2$
- B25: s.e.
- B28: critical t:
- B31: t-test:
- B34: Result:
- B39: Conclusion: (see Fig. 5.15)

Fig. 5.15 Wholesaler Price Comparison Formula Labels for Two-group t-test

Group	n	Mean	STDEV
1 Regular Price	7	117.7	19.9
2 Reduced Price	8	125.1	15.1
Null hypothesis:			
Research hypothesis:			
$(n_1 - 1) \times \text{STDEV1 squared}$			
$(n_2 - 1) \times \text{STDEV2 squared}$			
$n_1 + n_2 - 2$			
$1/n_1 + 1/n_2$			
s.e.			
critical t			
t-test			
Result:			
Conclusion:			

You now need to compute the values of the above formulas in the following cells:

E13: the result of the formula needed to compute cell B13 (use two decimals)

E16: the result of the formula needed to compute cell B16 (use two decimals)

E19: the result of the formula needed to compute cell B19

E22: the result of the formula needed to compute cell B22 (use two decimals)

E25: =SQRT(((E13+E16)/E19)*E22) (no spaces between)

Note the three open-parentheses after SQRT, and the three closed parentheses on the right side of this formula. You need three open parentheses and three closed parentheses in this formula or the formula will not work correctly.

The above formula gives a standard error of the difference of the means equal to 9.05 (two decimals).

E28: enter the critical t value from the t-table in Appendix E in this cell using $df = n_1 + n_2 - 2$ to find the critical t value

E31: =(D4-D5)/E25 (no spaces between)

Note that you need an open-parenthesis *before* D4 and a closed-parenthesis *after* D5 so that this answer of -7.40 is *THEN* divided by the standard error of the difference of the means of 9.05, to give a t-test value of -0.82 (note the minus sign here). Use two decimal places for the t-test result (see Fig. 5.16).

Fig. 5.16 Wholesaler Price Comparison Two-group t-test Formula Results

Group	n	Mean	STDEV
1 Regular Price	7	117.7	19.9
2 Reduced Price	8	125.1	15.1
Null hypothesis:			
Research hypothesis:			
(n1 - 1) x STDEV1 squared			2376.06
(n2 - 1) x STDEV2 squared			1596.07
n1 + n2 - 2			13
1/n1 + 1/n2			0.27
s.e.			9.05
critical t			2.160
t-test			-0.82
Result:			
Conclusion:			

Now write the following sentence in D34 to D37 to summarize the *result* of the study:

- D34: Since the absolute value
- D35: of t of -0.82 is less than
- D36: the critical t of 2.160, we
- D37: accept the null hypothesis.

Finally, write the following sentence in D39 to D43 to summarize the conclusion of the study:

- D39: There was no difference
- D40: in the number of units of

- D41: MP3 players sold at the
- D42: two prices. So, you should
- D43: not reduce the price!

Save your file as: MP4

Print the final spreadsheet so that it fits onto one page.

Write the null hypothesis and the research hypothesis by hand on your printout.

The final spreadsheet appears in Fig. 5.17.

Group	n	Mean	STDEV
1 Regular Price	7	117.7	19.9
2 Reduced Price	8	125.1	15.1
Null hypothesis:		$\mu_1 = \mu_2$	
Research hypothesis:		$\mu_1 \neq \mu_2$	
(n1 - 1) x STDEV1 squared			2376.06
(n2 - 1) x STDEV2 squared			1596.07
n1 + n2 - 2			13
1/n1 + 1/n2			0.27
s.e.			9.05
critical t			2.160
t-test			-0.82
Result:		Since the absolute value of t of - 0.82 is less than the critical t of 2.160, we accept the null hypothesis.	
Conclusion:		There was no difference in the number of units of MP3 players sold at the two prices. So, you should not reduce the price!	

Fig. 5.17 Wholesaler Price Comparison Final Spreadsheet

5.4 End-of-Chapter Practice Problems

- Suppose Boeing Company has hired you to do data analysis for its surveys that have been returned for its Morale Surveys that they had their managers answer during the past month. The items were summed to form a total score, in which a high score indicates high job satisfaction, while a low score indicates low job satisfaction. You select a random sample of managers, 202 females who averaged 84.80 on this survey with a standard deviation of 5.10. You also select a random sample of 241 males on this survey and they averaged 88.20 with a standard deviation of 4.30.
 - State the null hypothesis and the research hypothesis on an Excel spreadsheet.
 - Find the standard error of the difference between the means using Excel
 - Find the critical t value using Appendix E, and enter it on your spreadsheet.
 - Perform a t-test on these data using Excel. What is the value of t that you obtain?
 Use three decimal places for all figures in the formula section of your spreadsheet.
 - State your result on your spreadsheet.
 - State your conclusion in plain English on your spreadsheet.
 - Save the file as: Boeing3
- Suppose that you have been asked by the Director of the MS in Advertising program at the University of Illinois—Urbana to “run the data” to see if there is a gender difference in cumulative grade-point averages (GPAs) of MS in Advertising students who have completed all of the required courses for this degree. The Director has obtained the cooperation of the Registrar and has promised to keep the GPA information confidential. You want to test your Excel skills on some hypothetical data to make sure that you can do this analysis. These data appear in Fig. 5.18:

UNIVERSITY OF ILLINOIS -- URBANA	
GPA OF MS IN ADVERTISING STUDENTS WHO HAVE COMPLETED ALL ADVERTISING REQUIRED COURSES	
MALES	FEMALES
2.45	2.83
2.53	2.74
2.64	2.86
2.72	3.32
2.85	3.36
2.96	3.64
3.01	3.56
3.11	3.56
3.24	3.64
3.35	3.37
3.36	3.67
3.38	3.91
3.21	3.92
3.52	3.64
3.64	3.71
3.75	
3.86	

Fig. 5.18 Worksheet Data for Chap. 5: Practice Problem #2

- (a) State the null hypothesis and the research hypothesis on an Excel spreadsheet.
 - (b) Find the standard error of the difference between the means using Excel
 - (c) Find the critical t value using Appendix E, and enter it on your spreadsheet.
 - (d) Perform a t-test on these data using Excel. What is the value of t that you obtain?
 - (e) State your result on your spreadsheet.
 - (f) State your conclusion in plain English on your spreadsheet.
 - (g) Save the file as: GPA11
3. American Airlines offered an in-flight meal that passengers could purchase for \$11.00, and asked these customers to fill out a survey giving their opinion of the meal. Passengers were asked to rate their likelihood of purchasing this meal on a future flight on a 5-point scale. But, suppose that you have convinced the airline to change its survey item on purchase intention to a 7-point scale instead; the intention-to-buy item would then take the form in Fig. 5.19:

American Airlines Survey							
Item #10:	How likely are you to purchase an in-flight meal on a future flight?						
	1	2	3	4	5	6	7
	Definitely would not purchase						Definitely would purchase

Fig. 5.19 Worksheet Data for Chap. 5: Practice Problem #3

Passengers were asked on the survey to indicate whether they were either business travelers or vacationers. Suppose that the average rating last month for 64 “business travelers” was 3.23 with a standard deviation of 1.04, while the 56 “vacationers” had an average rating of 2.36 with a standard deviation of 1.35.

- (a) State the null hypothesis and the research hypothesis on an Excel spreadsheet.
- (b) Find the standard error of the difference between the means using Excel
- (c) Find the critical t value using Appendix E, and enter it on your spreadsheet.
- (d) Perform a t-test on these data using Excel. What is the value of t that you obtain?
- (e) State your result on your spreadsheet.
- (f) State your conclusion in plain English on your spreadsheet.
- (g) Save the file as: AAmeal3

References

- Keller, G. Statistics for Management and Economics (8th ed.). Mason, OH: South-Western Cengage Learning, 2009.
- Zikmund, W.G. and Babin, B.J. Exploring Marketing Research (10th ed.). Mason, OH: South-Western Cengage Learning, 2010.

Chapter 6

Correlation and Simple Linear Regression



There are many different types of “correlation coefficients,” but the one we will use in this book is the Pearson product-moment correlation which we will call: r .

6.1 What Is a “Correlation?”

Basically, a correlation is a number between -1 and $+1$ that summarizes the relationship between two variables, which we will call X and Y .

A correlation can be either positive or negative. *A positive correlation means that as X increases, Y increases. A negative correlation means that as X increases, Y decreases.* In statistics books, this part of the relationship is called the *direction* of the relationship (i.e., it is either positive or negative).

The correlation also tells us the *magnitude* of the relationship between X and Y . As the correlation approaches closer to $+1$, we say that the relationship is *strong and positive*.

As the correlation approaches closer to -1 , we say that the relationship is *strong and negative*.

A zero correlation means that there is no relationship between X and Y . This means that neither X nor Y can be used as a predictor of the other.

A good way to understand what a correlation means is to see a “picture” of the scatterplot of points produced in a chart by the data points. Let’s suppose that you want to know if variable X can be used to predict variable Y . We will place *the predictor variable X on the x -axis* (the horizontal axis of a chart) and *the criterion variable Y on the y -axis* (the vertical axis of a chart). Suppose, further, that you have collected data given in the scatterplots below (see Fig. 6.1 through Fig. 6.6).

Figure 6.1 shows the scatterplot for a perfect positive correlation of $r = +1.0$. This means that you can perfectly predict each y -value from each x -value because the data points move “upward-and-to-the-right” along a perfectly-fitting straight line (see Fig. 6.1)

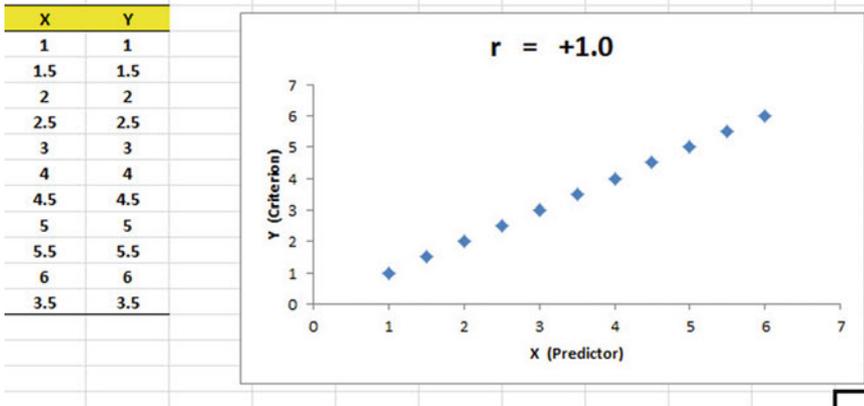


Fig. 6.1 Example of a Scatterplot for a Perfect, Positive Correlation ($r = +1.0$)

Figure 6.2 shows the scatterplot for a moderately positive correlation of $r = +.53$. This means that each x-value can predict each y-value moderately well because you can draw a picture of a “football” around the outside of the data points that move upward-and-to-the-right, but not along a straight line (see Fig. 6.2).

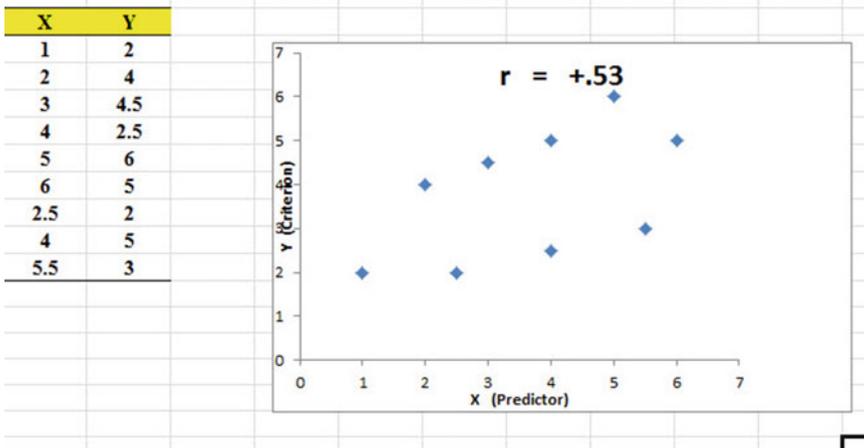


Fig. 6.2 Example of a Scatterplot for a Moderate, Positive Correlation ($r = +.53$)

Figure 6.3 shows the scatterplot for a low, positive correlation of $r = +.23$. This means that each x-value is a poor predictor of each y-value because the “picture” you could draw around the outside of the data points approaches a circle in shape (see Fig. 6.3).

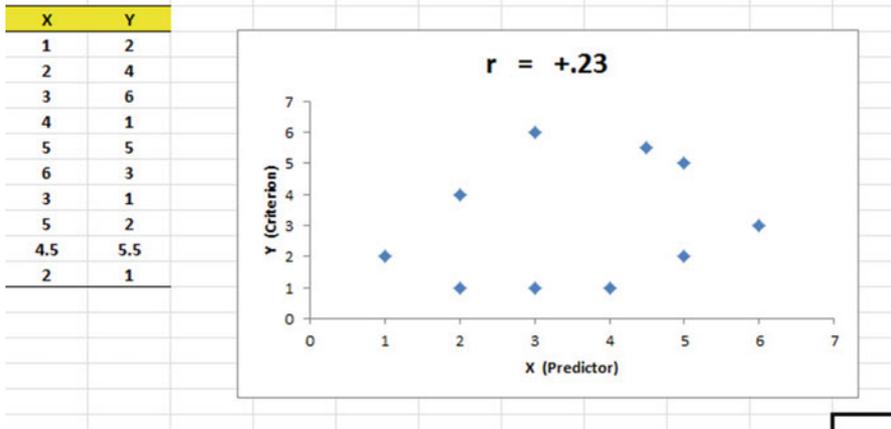


Fig. 6.3 Example of a Scatterplot for a Low, Positive Correlation ($r = +.23$)

We have not shown a Figure of a zero correlation because it is easy to imagine what it looks like as a scatterplot. A zero correlation of $r = .00$ means that there is no relationship between X and Y and the “picture” drawn around the data points would be a perfect circle in shape, indicating that you cannot use X to predict Y because these two variables are not correlated with one another.

Figure 6.4 shows the scatterplot for a low, negative correlation of $r = -.22$ which means that each X is a poor predictor of Y in an inverse relationship, meaning that as X increases, Y decreases (see Fig. 6.4). In this case, it is a negative correlation because the “football” you could draw around the data points slopes down and to the right.

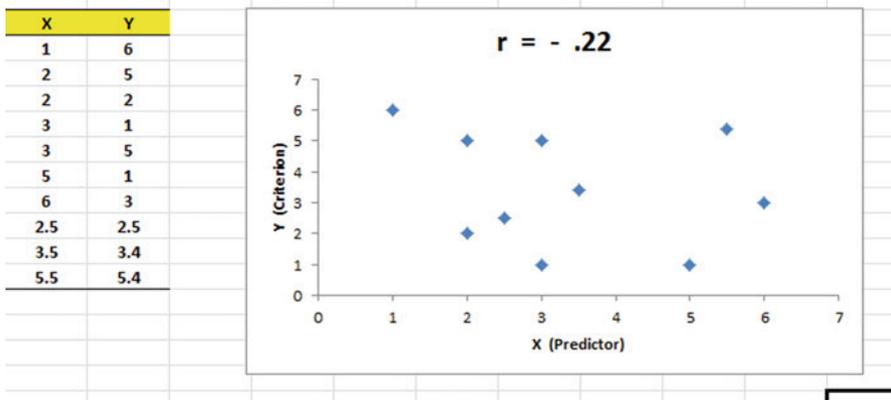


Fig. 6.4 Example of a Scatterplot for a Low, Negative Correlation ($r = -.22$)

Figure 6.5 shows the scatterplot for a moderate, negative correlation of $r = -.39$ which means that X is a moderately good predictor of Y, although there is an inverse relationship between X and Y (i.e., as X increases, Y decreases; see Fig. 6.5). In this case, it is a negative correlation because the “football” you could draw around the data points slopes down and to the right.

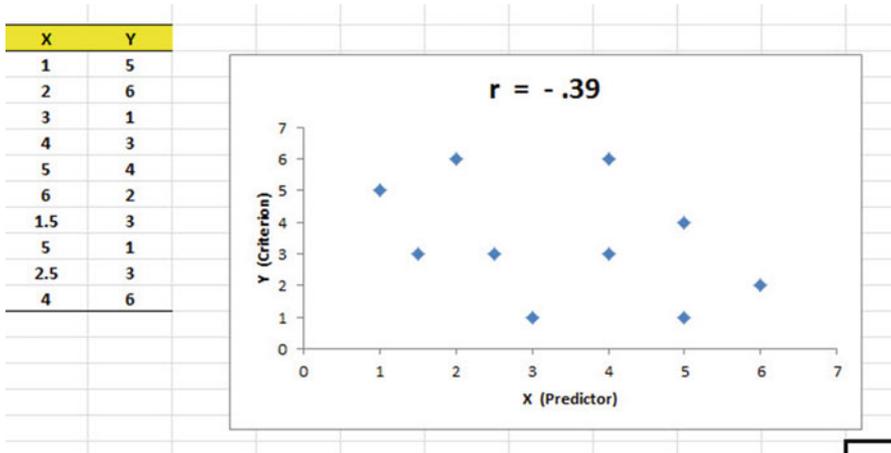


Fig. 6.5 Example of a Scatterplot for a Moderate, Negative Correlation ($r = -.39$)

Figure 6.6 shows a perfect negative correlation of $r = -1.0$ which means that X is a perfect predictor of Y, although in an inverse relationship such that as X increases, Y decreases. The data points fit perfectly along a downward-sloping straight line (see Fig. 6.6).

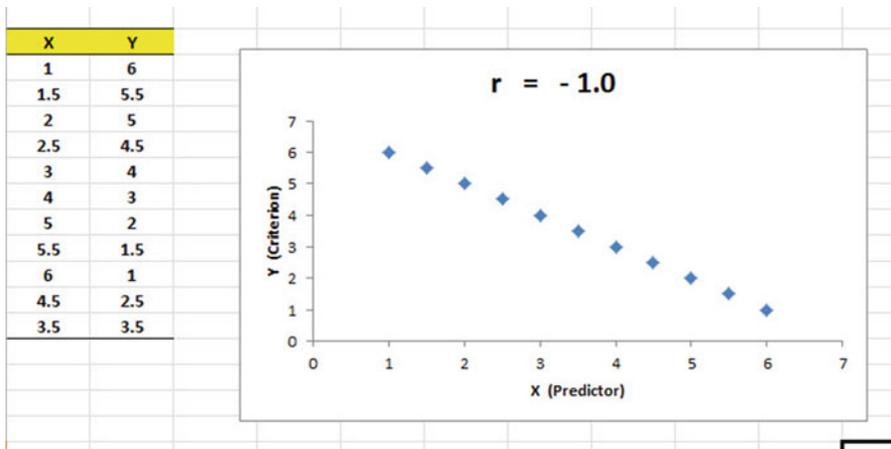


Fig. 6.6 Example of a Scatterplot for a Perfect, Negative Correlation ($r = -1.0$)

Let’s explain the formula for computing the correlation r so that you can understand where the number summarizing the correlation came from.

In order to help you to understand *where* the correlation number that ranges from -1.0 to $+1.0$ comes from, we will walk you through the steps involved to use the formula as if you were using a pocket calculator. This is the one time in this book that we will ask you to use your pocket calculator to find a correlation, but knowing how the correlation is computed step-by-step will give you the opportunity to understand *how* the formula works in practice.

To do that, let’s create a situation in which you need to find the correlation between two variables.

Suppose that you have been hired by a manager of a supermarket chain to find the relationship between the amount of money spent weekly by the chain on television ads and the weekly sales of the supermarket chain in St. Louis. You collect the data from the past 8 weeks given in Fig. 6.7.

	Week	TV ad cost (\$000)	Weekly Sales (\$000)
	1	4.8	94
	2	1.9	87
	3	3.8	93
	4	2.3	89
	5	2.9	92
	6	3.3	92
	7	2.4	93
	8	2.8	92
n		8	8
MEAN		3.03	91.50
STDEV		0.93	2.33

Fig. 6.7 Worksheet Data for a Supermarket Chain (Practical Example)

For the purposes of explanation, let’s call the weekly cost of TV ads as the predictor variable X , and the weekly sales as the criterion variable Y . Notice that the data for the cost of TV ads for each week is in thousands of dollars (\$000). For example, the TV ads for week 6 cost \$3300, and when we “move the decimal place three places to the left to change the amount to thousands of dollars,” this becomes 3.3. Similarly, the weekly sales for week 6 were really \$92,000 as those data are also in thousands of dollars format (\$000).

Notice also that we have used Excel to find the sample size for both variables, X and Y , and the MEAN and STDEV of both variables. (You can practice your Excel skills by seeing if you get these same results when you create an Excel spreadsheet for these data.)

Now, let's use the above table to compute the correlation r between the weekly cost of TV ads and the weekly sales of this supermarket chain using your pocket calculator.

6.1.1 Understanding the Formula for Computing a Correlation

Objective: To understand the formula for computing the correlation r

The formula for computing the correlation r is as follows:

$$r = \frac{\frac{1}{n-1} \sum (X - \bar{X})(Y - \bar{Y})}{S_x S_y} \quad (6.1)$$

This formula looks daunting at first glance, but let's "break it down into its steps" to understand how to compute the correlation r .

6.1.2 Understanding the Nine Steps for Computing a Correlation, r

Objective: To understand the nine steps of computing a correlation r

The nine steps are as follows:

Step	Computation	Result
1	Find the sample size n by noting the number of weeks	8
2	Divide the number 1 by the sample size minus 1 (i.e., $1/7$)	0.14286
3	<i>For each week</i> , take the cost of TV ads for that week and subtract the mean cost of TV ads for the 8 weeks and call this $X - \bar{X}$ (For example, for week 6, this would be: $3.3 - 3.03$)	0.27
	Note: With your calculator, this difference is 0.27, but when Excel uses 16 decimal places for every computation, this result will be 0.28 instead of 0.27.	
4	<i>For each week</i> , take the weekly sales for that week and subtract the mean weekly sales for the 8 weeks and call this $Y - \bar{Y}$ (For example, for week 6, this would be: $92 - 91.50$)	0.50
5	Then, <i>for each week</i> , multiply $(X - \bar{X})$ times $(Y - \bar{Y})$ (For example, for week 6 this would be: 0.27×0.50)	0.135
6	Add the results of $(X - \bar{X})$ times $(Y - \bar{Y})$ for the 8 weeks	11.50

	X		Y		
Week	TV ad cost (\$000)	Weekly Sales (\$000)	$X - \bar{X}$	$Y - \bar{Y}$	$(X - \bar{X})(Y - \bar{Y})$
1	4.8	94	1.78	2.50	4.44
2	1.9	87	-1.13	-4.50	5.06
3	3.8	93	0.78	1.50	1.16
4	2.3	89	-0.73	-2.50	1.81
5	2.9	92	-0.13	0.50	-0.06
6	3.3	92	0.28	0.50	0.14
7	2.4	93	-0.63	1.50	-0.94
8	2.8	92	-0.23	0.50	-0.11
n	8	8		Total	11.50
MEAN	3.03	91.50			
STDEV	0.93	2.33			

Fig. 6.8 Worksheet for Computing the Correlation, r

Steps 1–6 would produce the Excel table given in Fig. 6.8.

Notice that when Excel multiplies a minus number by a minus number, the result is a plus number (for example for week 2: $(-1.13 \times -4.50 = +5.06)$. And when Excel multiplies a minus number by a plus number, the result is a negative number (for example for week 5: $(-0.13 \times +0.50 = -0.06)$.

Note: Excel computes all computation to 16 decimal places. So, when you check your work with a calculator, you frequently get a slightly different answer than Excel’s answer.

For example, when you compute above:

$$(X - \bar{X}) \times (Y - \bar{Y}) \text{ for Week 2, your calculator gives:}$$

$$(-1.13) \times (-4.50) = +5.085, \tag{6.2}$$

But, as you can see from the table, Excel’s answer of 5.06 is *more accurate* because Excel uses 16 decimal places for every number.

You should also note that when you do Step 6, you have to be careful to add all of the positive numbers first to get +12.61 and then add all of the negative numbers second to get -1.11, so that when you subtract these two numbers you get +11.50 as your answer to Step 6.

Step		
7	Multiply the answer for step 2 above by the answer for step 6 (0.14286 × 11.5)	1.6429
8	Multiply the STDEV of X times the STDEV of Y (0.93 × 2.33)	2.1669
9	Finally, divide the answer from step 7 by the answer from step 8 (1.6429 divided by 2.1669)	+0.76

This number of 0.76 is the correlation between the weekly cost of TV ads (X) and the weekly sales in this supermarket chain (Y) over this 8-week period. The number $+0.76$ means that there is a strong, positive correlation between these two variables. That is, as the chain increases its spending on TV ads, its sales for that week increase. For a more detailed discussion of correlation, see Zikmund and Babin (2010).

You could also use the results of the above table in the formula for computing the correlation r in the following way:

$$\text{correlation } r = [1/(n - 1) \times \sum (X - \bar{X})(Y - \bar{Y})] / (\text{STDEV}_x \times \text{STDEV}_y)$$

$$\text{correlation } r = [(1/7) \times 11.50] / [(0.93) \times (2.33)]$$

$$\text{correlation } = r = 0.76$$

Now, let's discuss how you can use Excel to find the correlation between two variables in a much simpler, and much faster, fashion than using your calculator.

6.2 Using Excel to Compute a Correlation Between Two Variables

Objective: To use Excel to find the correlation between two variables

Suppose that you have been hired by the owner of a supermarket chain in St. Louis to make a recommendation as to how many shelf facings of Kellogg's Corn Flakes this chain should use. A "shelf facing" is the number of boxes of the cereal that are stacked beside one another. Thus a shelf facing of 3 means that three boxes of Kellogg's Corn Flakes are stacked beside each other on the supermarket shelf in the cereals section.

You randomly assign supermarket locations to your study, and you randomly select the number of facings used in each supermarket location, where the number of facings range from 1 to 3 facings. You track the weekly sales (in thousands of dollars) of this cereal over a 10-week period, and the resulting sales figures are given in Fig. 6.9.

Fig. 6.9 Worksheet Data for the Number of Facings and Sales (Practical Example)

Week	No. of facings	Sales (\$000)
1	1	1.1
2	2	2.2
3	3	2.1
4	1	1.2
5	2	2.3
6	3	5.2
7	3	4.6
8	2	2.3
9	2	1.9
10	3	4.5

You want to determine if there is a *relationship* between the number of facings of Kellogg’s Corn Flakes and the weekly sales of this cereal, and you decide to use a correlation to determine this relationship. Let’s call the number of facings, X , and the sales figures, Y .

Create an Excel spreadsheet with the following information:

A2: Week

B2: No. of facings

C2: Sales (\$000)

A3: 1

Next, change the width of Columns B and C so that the information fits inside the cells.

Now, complete the remaining figures in the table given above so that A12 is 10, B12 is 3, and C12 is 4.5 (Be sure to double-check your figures to make sure that they are correct!) Then, center the information in all of these cells.

A14: n

A15: mean

A16: stdev

Next, define the “name” to the range of data from B3:B12 as: facings

We discussed earlier in this book (see Sect. 1.4.4) how to “name a range of data,” but here is a reminder of how to do that:

To give a “name” to a range of data:

Click on the top number in the range of data and drag the mouse down to the bottom number of the range.

For example, to give the name: “facings” to the cells: B3:B12, click on B3, and drag the pointer down to B12 so that the cells B3:B12 are highlighted on your computer screen. Then, click on:

Formulas

Define name (top center of your screen)

facings (in the Name box; see Fig. 6.10)

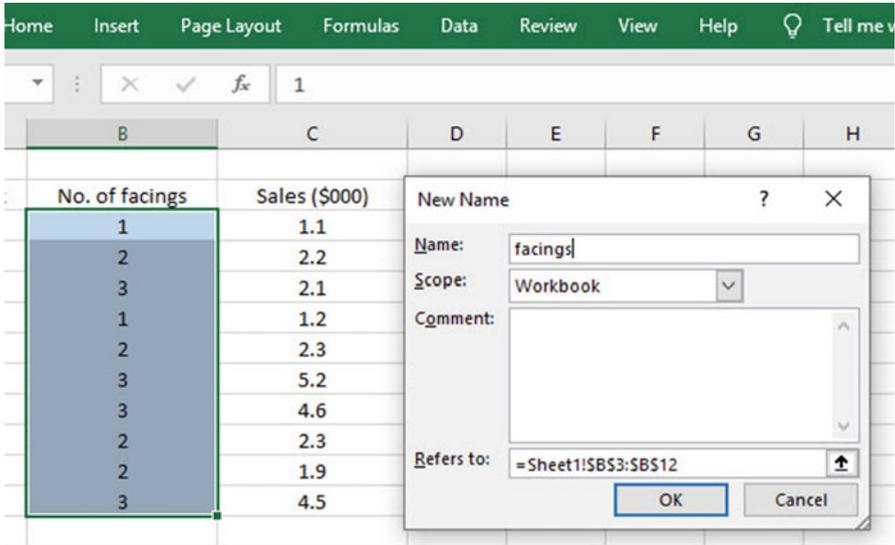


Fig. 6.10 Dialogue Box for Naming a Range of Data as: “facings”

OK

Now, repeat these steps to give the name: *sales* to C3:C12

Finally, click on any blank cell on your spreadsheet to “deselect” cells C3:C12 on your computer screen.

Now, complete the data for these sample sizes, means, and standard deviations in columns B and C so that B16 is 0.79, and C16 is 1.47 (use two decimals for the means and standard deviations; center the data in B14:C16; see Fig. 6.11)

Fig. 6.11 Example of Using Excel to Find the Sample Size, Mean, and STDEV

Week	No. of facings	Sales (\$000)
1	1	1.1
2	2	2.2
3	3	2.1
4	1	1.2
5	2	2.3
6	3	5.2
7	3	4.6
8	2	2.3
9	2	1.9
10	3	4.5
n	10	10
mean	2.20	2.74
stdev	0.79	1.47

Objective: Find the correlation between the number of facings and the weekly sales dollars.

B18: correlation

C18: =correl(facings,sales) ; see Fig. 6.12

	A	B	C	D	E
1					
2	Week	No. of facings	Sales (\$000)		
3	1	1	1.1		
4	2	2	2.2		
5	3	3	2.1		
6	4	1	1.2		
7	5	2	2.3		
8	6	3	5.2		
9	7	3	4.6		
10	8	2	2.3		
11	9	2	1.9		
12	10	3	4.5		
13					
14	n	10	10		
15	mean	2.20	2.74		
16	stdev	0.79	1.47		
17					
18		correlation	=correl(facings,sales)		
19					

Fig. 6.12 Example of Using Excel's =correl Function to Compute the Correlation Coefficient

Hit the Enter key to compute the correlation

C18: format this cell to two decimals

Note that the equal sign tells Excel that you are going to use a formula.

The correlation between the number of facings (X) and weekly sales (Y) is $+0.83$, a very strong positive correlation. This means that you have evidence that there is a strong relationship between these two variables. In effect, the more facings (when 1, 2, 3 facings are used), the higher the weekly sales dollars generated for this cereal.

Save this file as: FACINGS5

The final spreadsheet appears in Fig. 6.13.

Fig. 6.13 Final Result of Using the =correl Function to Compute the Correlation Coefficient

C18		fx =CORREL(facings,sales)		
A	B	C	D	E
Week	No. of facings	Sales (\$000)		
1	1	1.1		
2	2	2.2		
3	3	2.1		
4	1	1.2		
5	2	2.3		
6	3	5.2		
7	3	4.6		
8	2	2.3		
9	2	1.9		
10	3	4.5		
n	10	10		
mean	2.20	2.74		
stdev	0.79	1.47		
	correlation	0.83		

6.3 Creating a Chart and Drawing the Regression Line onto the Chart

This section deals with the concept of “linear regression.” Technically, the use of a simple linear regression model (i.e., the word “simple” means that only one predictor, X, is used to predict the criterion, Y) requires that the data meet the following four assumptions if that statistical model is to be used:

1. The underlying relationship between the two variables under study (X and Y) is *linear* in the sense that a straight line, and not a curved line, can fit among the data points on the chart.
2. The errors of measurement are independent of each other (e.g. the errors from a specific time period are sometimes correlated with the errors in a previous time period).
3. The errors fit a normal distribution of Y-values at each of the X-values.
4. The variance of the errors is the same for all X-values (i.e., the variability of the Y-values is the same for both low and high values of X).

A detailed explanation of these assumptions is beyond the scope of this book, but the interested reader can find a detailed discussion of these assumptions in Levine et al. (2011, pp. 529–530).

Now, let's create a chart summarizing these data.

Important note: Whenever you draw a chart, it is ESSENTIAL that you put the predictor variable (X) on the left, and the criterion variable (Y) on the right in your Excel spreadsheet, so that you know which variable is the predictor variable and which variable is the criterion variable. If you do this, you will save yourself a lot of grief whenever you do a problem involving correlation and simple linear regression using Excel!

Important note: You need to understand that in any chart that has one predictor and a criterion that there are really TWO LINES that can be drawn between the data points:

- (1) One line uses X as the predictor, and Y as the criterion
- (2) A second line uses Y as the predictor, and X as the criterion

This means that you have to be very careful to note in your input data the cells that contain X as the predictor, and Y as the criterion. If you get these cells mixed up and reverse them, you will create the wrong line for your data and you will have botched the problem terribly.

This is why we STRONGLY RECOMMEND IN THIS BOOK that you always put the X data (i.e., the predictor variable) on the LEFT of your table, and the Y data (i.e., the criterion variable) on the RIGHT of your table on your spreadsheet so that you don't get these variables mixed up.

Also note that the correlation, r , will be exactly the same correlation no matter which variable you call the predictor variable and which variable you call the criterion variable. The correlation coefficient just summarizes the relationship between two variables, and doesn't care which one is the predictor and which one is the criterion.

Let's suppose that you would like to use the number of facings of Corn Flakes as the predictor variable, and that you would like to use it to predict the weekly sales dollars of this cereal. Since the correlation between these two variables is $+0.83$, this shows that there is a strong, positive relationship and that the number of facings is a good predictor of the weekly sales for this cereal.

1. Open the file that you saved earlier in this chapter: FACINGS5

6.3.1 Using Excel to Create a Chart and the Regression Line Through the Data Points

Objective: To create a chart and the regression line summarizing the relationship between the number of shelf facings and the weekly sales (\$000).

2. Click and drag the mouse to highlight both columns of numbers (B3:C12), but do not highlight the labels at the top of Column B and Column C.

Highlight the data set: B3:C12

Insert (top left of screen)

Highlight: Scatter chart icon (immediately above the word: “Charts” at the top center of your screen)

Click on the down arrow on the right of the chart icon

Highlight the top left scatter chart icon (see Fig. 6.14)

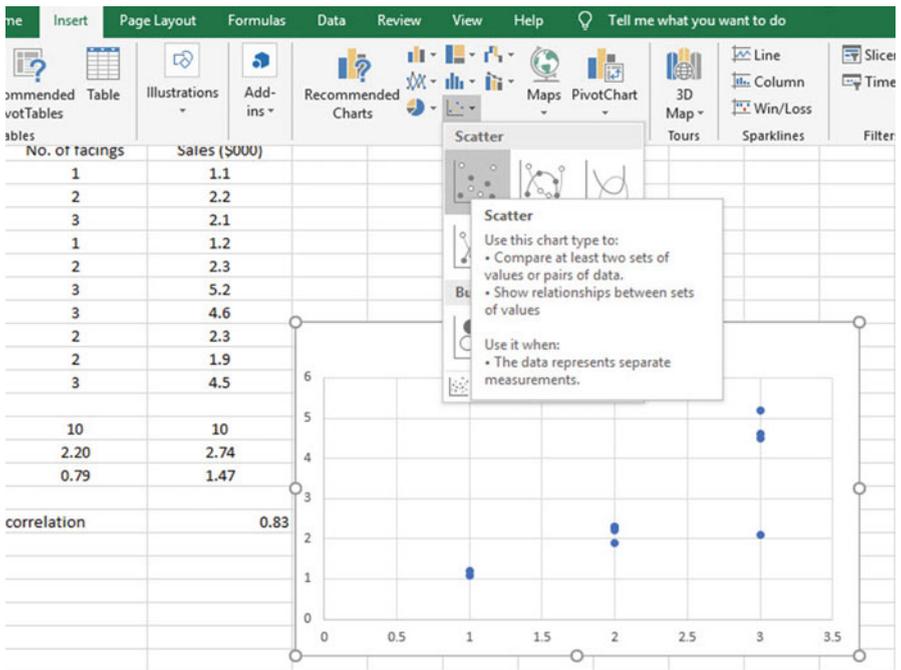


Fig. 6.14 Example of Selecting a Scatter Chart

Click on the top left chart to select it

Click on the “+ icon” to the right of the chart (CHART ELEMENTS).

Click on the check mark next to “Chart Title” **and also** next to “Gridlines” to remove these check marks (see Fig. 6.15)

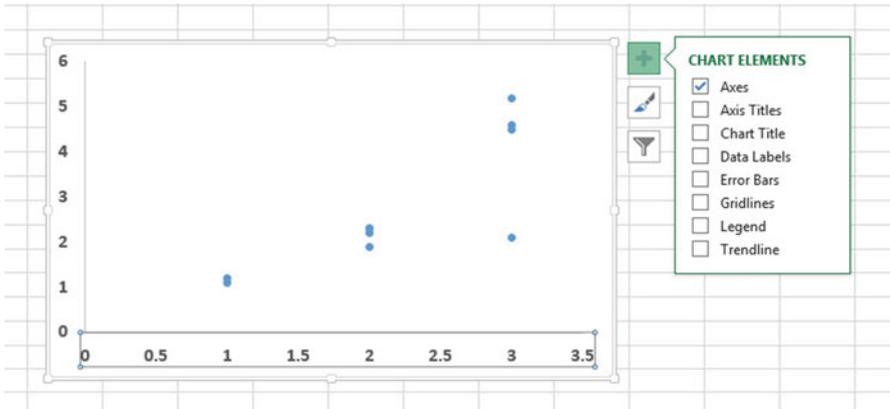


Fig. 6.15 Example of Chart Elements Selected

Click on the box next to: “Chart Title” and then click on the arrow to its right. Then, click on: “Above chart”.

Note that the words: “Chart Title” are now in a box at the top of the chart (see Fig. 6.16)

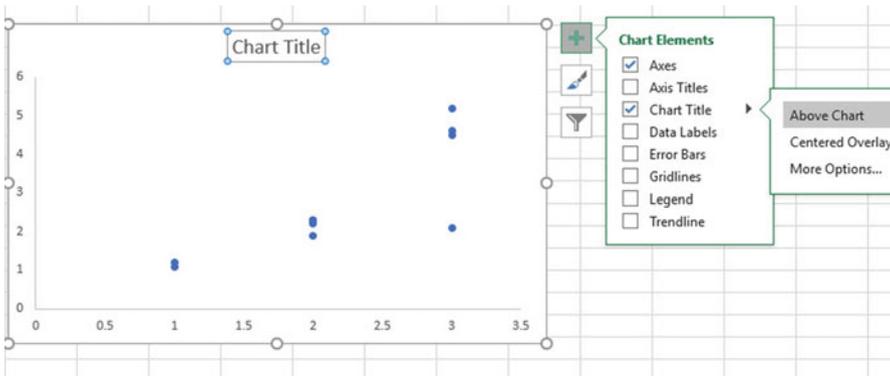


Fig. 6.16 Example of Chart Title Selected

Enter the following Chart Title to the right of f_x at the top of your screen: RELATIONSHIP BETWEEN NO. OF FACINGS AND SALES (see Fig. 6.17)

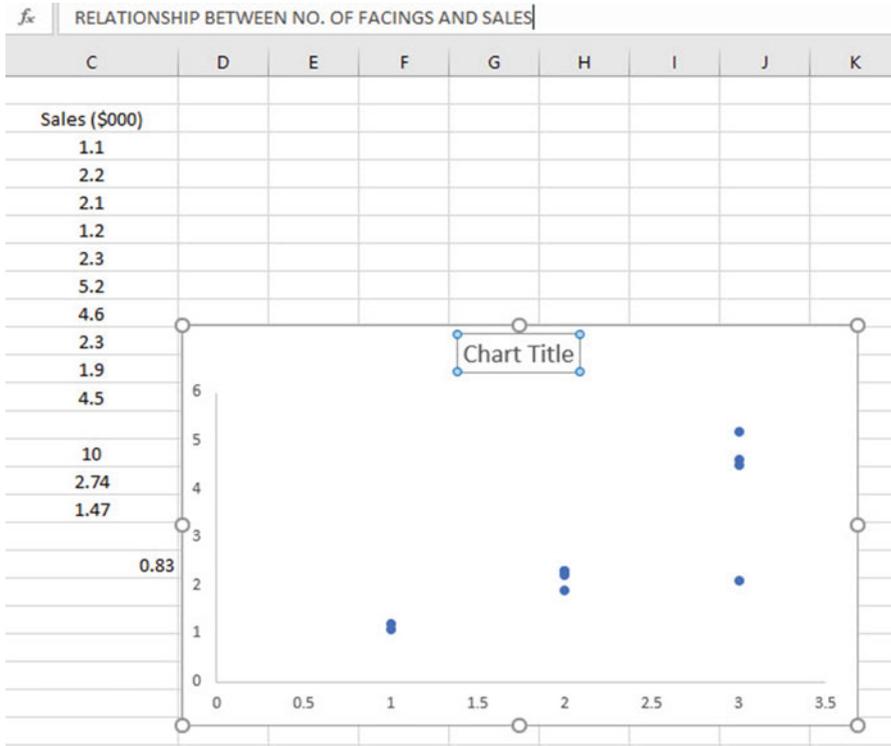


Fig. 6.17 Example of Creating a Chart Title

Hit the Enter Key to enter this chart title onto the chart

Click *inside the chart at the top right corner of the chart* to “deselect” the box around the Chart Title (see Fig. 6.18)

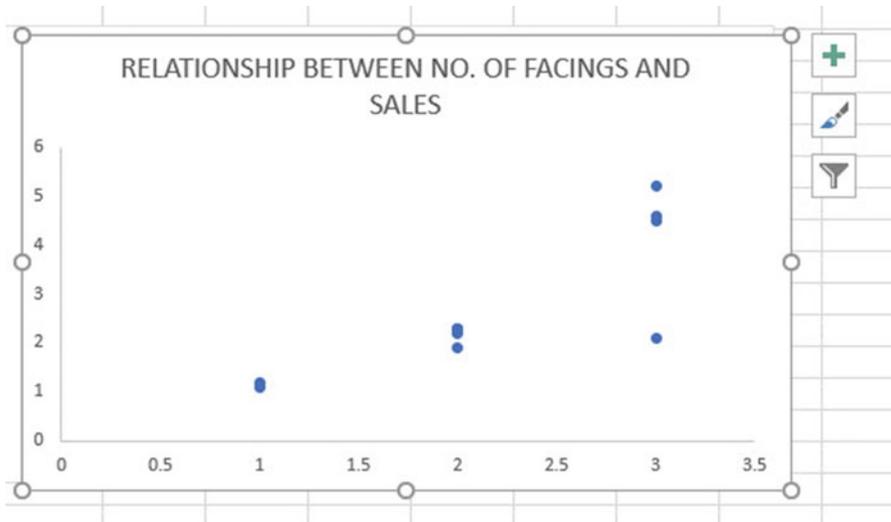


Fig. 6.18 Example of a Chart Title Inserted onto the Chart

Click on the “+ box” to the right of the chart

Add a check mark to the left of “Axis Titles” (This will create an “Axis Title” box on the y-axis of the chart)

Click on the right arrow for: “Axis titles” and then click on: “Primary Horizontal” to remove the check mark in its box (this will create the y-axis title)

Enter the following y-axis title to the right of f_x at the top of your screen:

SALES (\$000)

Then, hit the Enter Key to enter this y-axis title to the chart

Click *inside the chart at the top right corner of the chart* to “deselect” the box around the y-axis title (see Fig. 6.19)

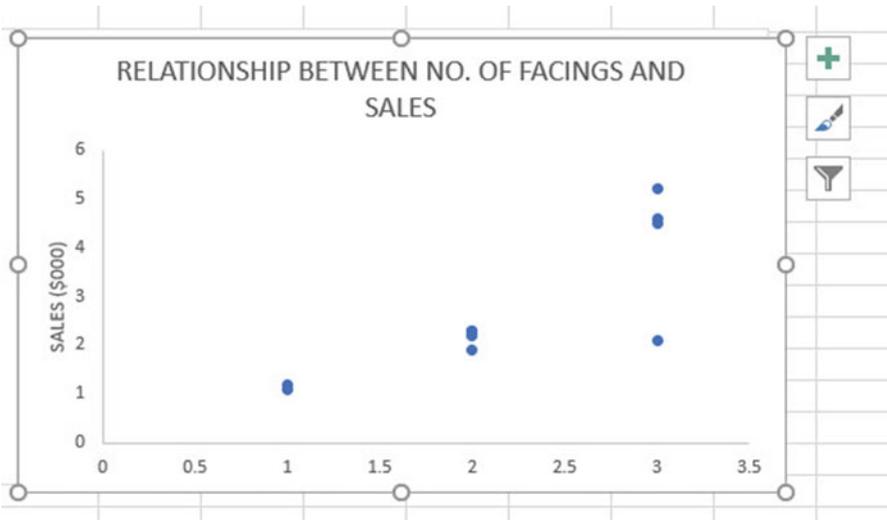


Fig. 6.19 Example of Adding a y-axis Title to the Chart

Click on the “+ box” to the right of the chart

Highlight: “Axis Titles” and click on its right arrow

Click on the words: “Primary Horizontal” to add a check mark to its box (this creates an “Axis Title” box on the x-axis of the chart)

Enter the following x-axis title to the right of f_x at the top of your screen:

NO. OF FACINGS

Then, hit the Enter Key to add this x-axis title to the chart

Click *inside the chart at the top right corner of the chart* to “deselect” the box around the x-axis title (see Fig. 6.20).

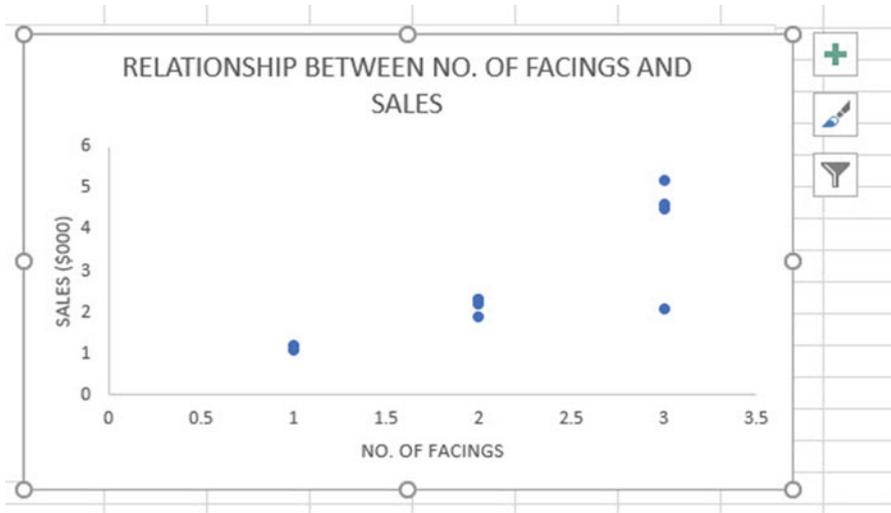


Fig. 6.20 Example of a Chart Title, an x-axis Title, and a y-axis Title

6.3.1.1 Drawing the Regression Line Through the Data Points in the Chart

Objective: To draw the regression line through the data points on the chart

Right-click on any one of the data points inside the chart

Highlight: Add Trendline (see Fig. 6.21)

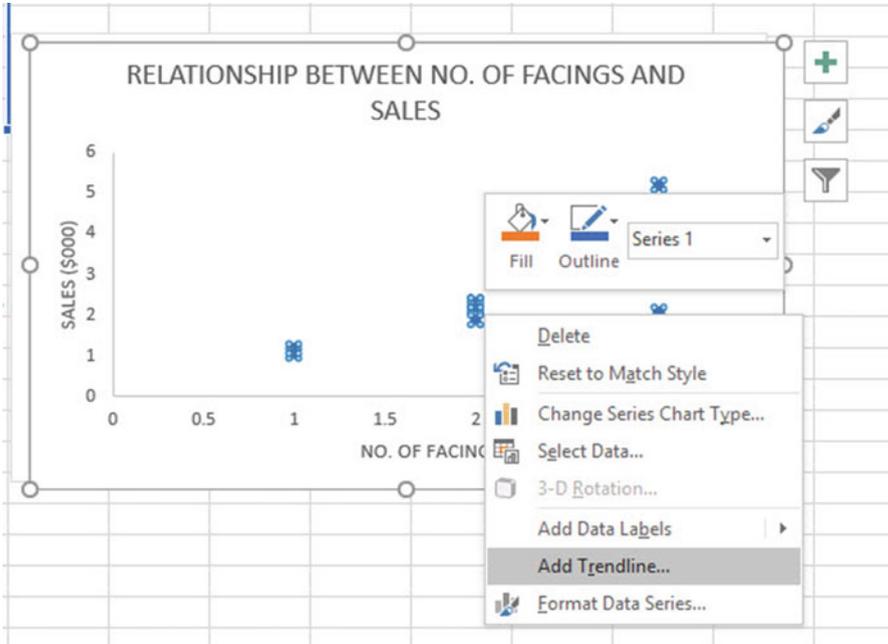


Fig. 6.21 Dialogue Box for Adding a Trendline to the Chart

Click on: Add Trendline

Linear (be sure the “linear” button near the top is selected on the “Format Trendline” dialog box; see Fig. 6.22)

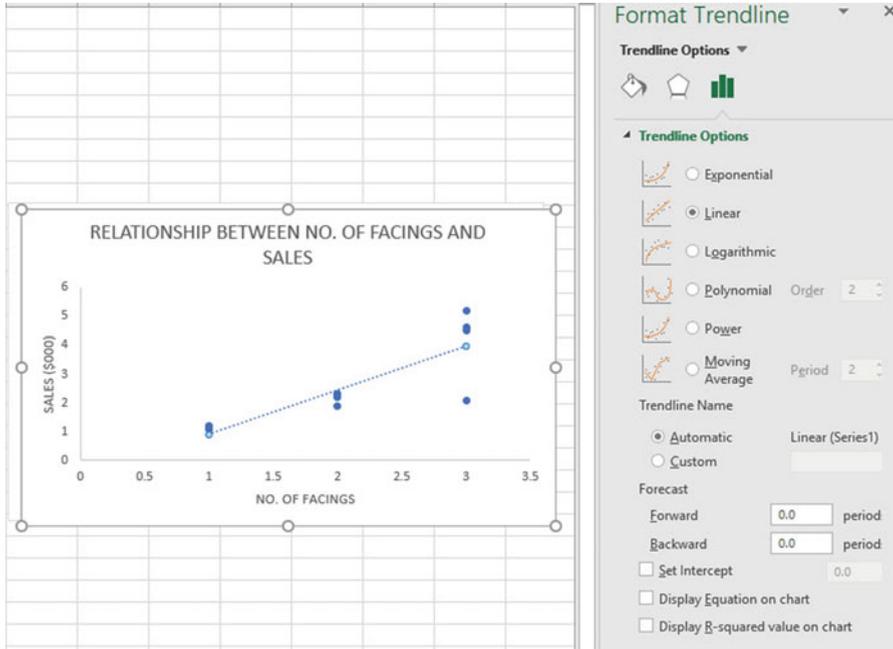


Fig. 6.22 Dialogue Box for a Linear Trendline

Click on the X at the top right of the “Format Trendline” dialog box to close this dialog box

Click on any blank cell *outside the chart* to “deselect” the chart

Save this file as: FACINGS7

Your spreadsheet should look like the spreadsheet in Fig. 6.23.

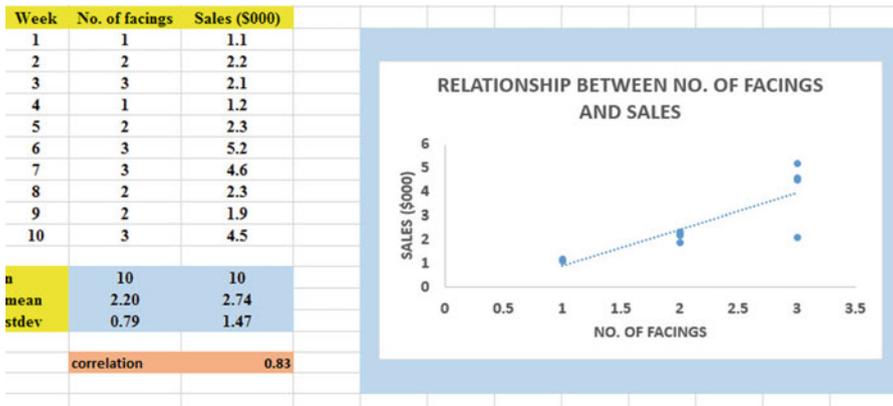


Fig. 6.23 Final Chart with the Trendline Fitted Through the Data Points of the Scatterplot

6.3.1.2 Moving the Chart Below the Table in the Spreadsheet

Objective: To move the chart below the table

Left-click your mouse on *any white space to the right of the top title inside the chart*, keep the left-click down, and drag the chart down and to the left so that the top left corner of the chart is in cell A20, then take your finger off the left-click of the mouse (see Fig. 6.24).

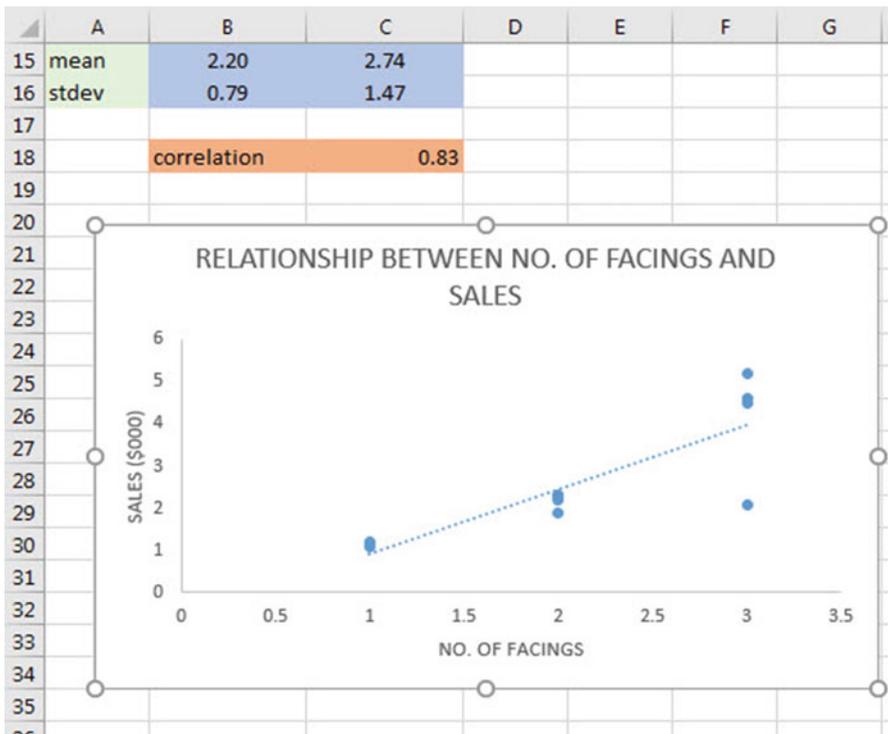


Fig. 6.24 Example of Moving the Chart Below the Table

6.3.1.3 Making the Chart “Longer” So That It Is “Taller”

Objective: To make the chart “longer” so that it is taller

Left-click your mouse on the bottom-center of the chart to create an “up-and-down-arrow” sign, hold the left-click of the mouse down and drag the bottom of the chart down to row 42 to make the chart longer, and then take your finger off the mouse.

6.3.1.4 Making the Chart “Wider”

Objective: To make the chart “wider”

Put the pointer at the middle of the right-border of the chart to create a “left-to-right arrow” sign, and then left-click your mouse and hold the left-click down while you drag the right border of the chart to the middle of Column H to make the chart wider.

Now, click on any blank cell outside the chart to “deselect” the chart (see Fig. 6.25).

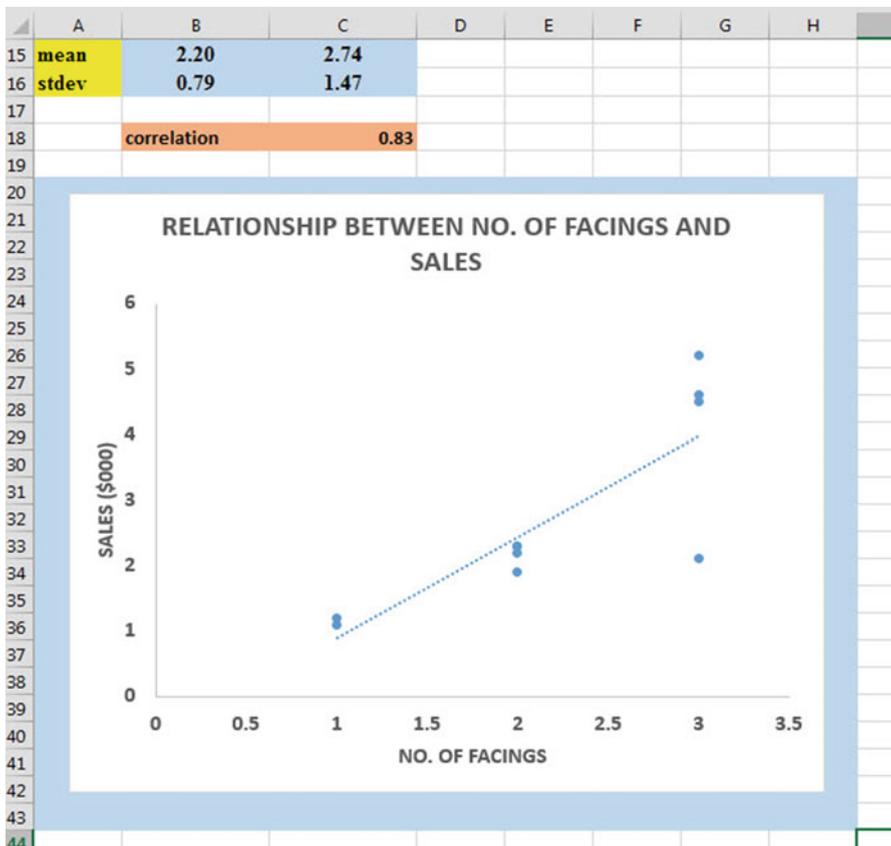


Fig. 6.25 Example of a Chart that is Enlarged to Fit the Cells: A20:H42

6.4 Printing a Spreadsheet So That the Table and Chart Fit onto One Page

Objective: To print the spreadsheet so that the table and the chart fit onto one page

Page Layout (top of screen)

Change the scale at the middle icon near the top of the screen “Scale to Fit” by clicking on the down-arrow until it reads “95%” so that the table and the chart will fit onto one page on your screen (see Fig. 6.26)

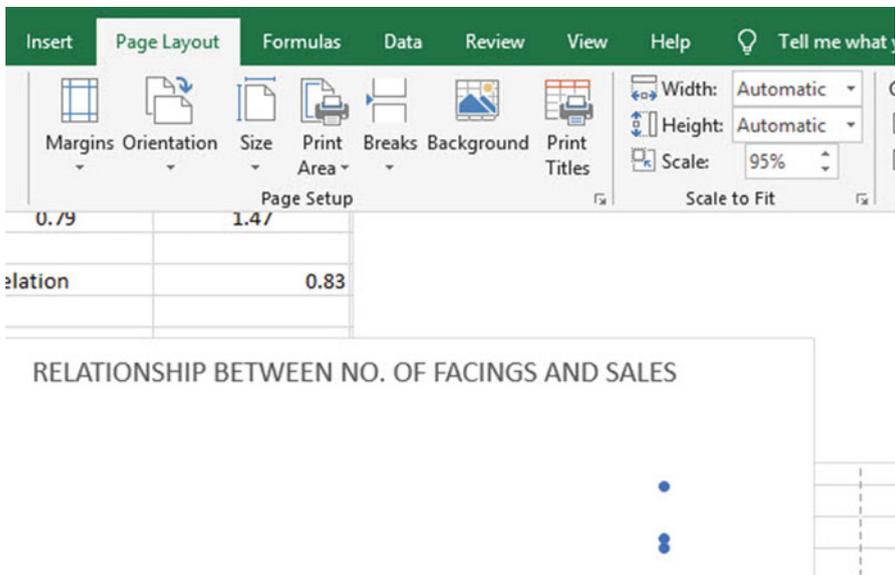


Fig. 6.26 Example of the Page Layout for Reducing the Scale of the Chart to 95% of Normal Size

File

Print

Print (see Fig. 6.27)

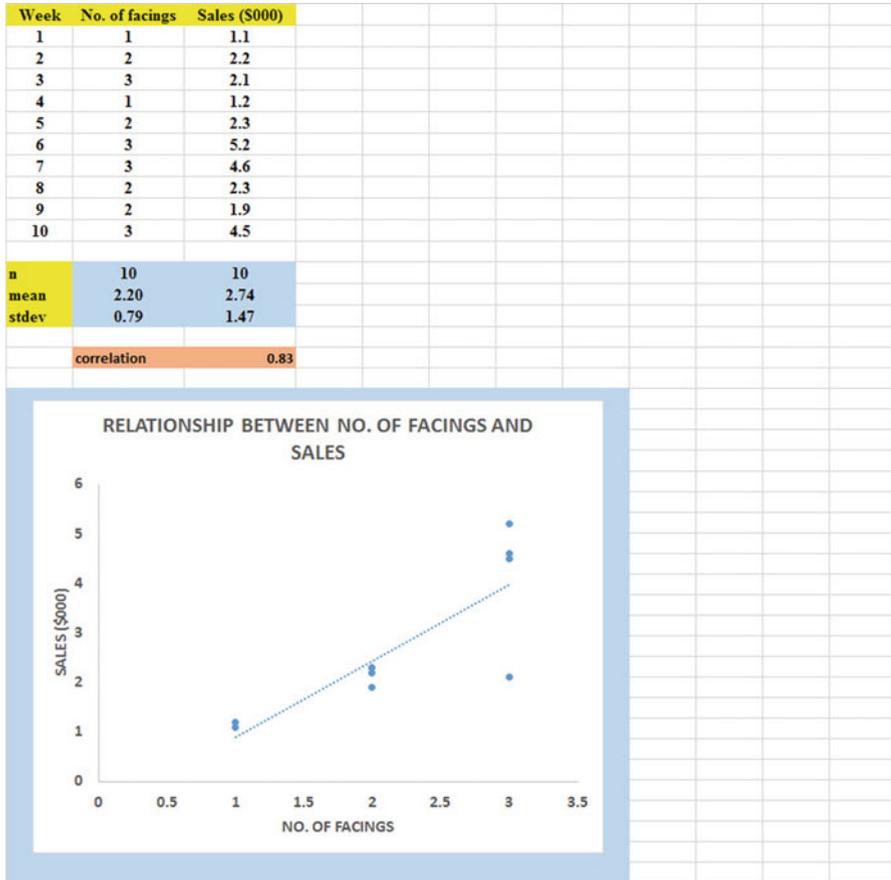


Fig. 6.27 Final Spreadsheet of a Table and a Chart (95% Scale to Fit Size)

Save your file as: FACINGS8

6.5 Finding the Regression Equation

The main reason for charting the relationship between X and Y (i.e., No. of facings as X and Sales (\$000) as Y in our example) is to see if there is a strong relationship between X and Y so that the regression equation that summarizes this relationship can be used to predict Y for a given value of X.

Since we know that the correlation between the number of facings and sales is $+0.83$, this tells us that it makes sense to use the number of facings to predict the weekly sales that we can expect based on past data.

We now need to find that regression equation that is the equation of the “best-fitting straight line” through the data points.

Objective: To find the regression equation summarizing the relationship between X and Y.

In order to find this equation, we need to check to see if your version of Excel contains the “Data Analysis ToolPak” necessary to run a regression analysis.

6.5.1 Installing the Data Analysis ToolPak into Excel

Objective: To install the Data Analysis ToolPak into Excel

Since there are currently several versions of Excel in the marketplace (2013, 2016, 2019), we will give a brief explanation of how to install the Data Analysis ToolPak into each of these versions of Excel.

6.5.1.1 Installing the Data Analysis ToolPak into Excel 2019

Open a new Excel spreadsheet

Click on: Data (at the top of your screen)

Look at the top of your monitor screen. Do you see the words: “Data Analysis” at the far right of the screen? If you do, the Data Analysis ToolPak for Excel 2019 was correctly installed when you installed Office 2019, and you should skip ahead to Sect. 6.5.2.

If the words: “Data Analysis” are not at the top right of your monitor screen, then the ToolPak component of Excel 2019 was not installed when you installed Office 2019 onto your computer. If this happens, you need to follow these steps:

File

Options (bottom left of screen)

Note: This creates a dialog box with “Excel Options” at the top left of the box

Add-Ins (on left of screen)

Manage: Excel Add-Ins (at the bottom of the dialog box)

Go (at bottom center of dialog box)

Highlight: Analysis ToolPak (in the Add-Ins dialog box)

Put a check mark to the left of Analysis Toolpak

OK (at the right of this dialog box)

Data

You now should have the words: “Data Analysis” at the top right of your screen to show that this feature has been installed correctly

Note: If these steps do not work, you should try these steps instead: File/Options (bottom left)/Add-ins/Analysis ToolPak/Go/click to the left of Analysis ToolPak to add a check mark/OK

If you need help doing this, ask your favorite “computer techie” for help.
You are now ready to skip ahead to Sect. [6.5.2](#)

6.5.1.2 Installing the Data Analysis ToolPak into Excel 2016

Open a new Excel spreadsheet

Click on: Data (at the top of your screen)

Look at the top of your monitor screen. Do you see the words: “Data Analysis” at the far right of the screen? If you do, the Data Analysis ToolPak for Excel 2016 was correctly installed when you installed Office 2016, and you should skip ahead to Sect. [6.5.2](#).

If the words: “Data Analysis” are not at the top right of your monitor screen, then the ToolPak component of Excel 2016 was not installed when you installed Office 2016 onto your computer. If this happens, you need to follow these steps:

File

Options (bottom left of screen)

Note: This creates a dialog box with “Excel Options” at the top left of the box

Add-Ins (on left of screen)

Manage: Excel Add-Ins (at the bottom of the dialog box)

Go (at bottom center of dialog box)

Highlight: Analysis ToolPak (in the Add-Ins dialog box)

Put a check mark to the left of Analysis Toolpak

OK (at the right of this dialog box)

Data

You now should have the words: “Data Analysis” at the top right of your screen to show that this feature has been installed correctly

Note: If these steps do not work, you should try these steps instead: File/Options (bottom left)/Add-ins/Analysis ToolPak/Go/click to the left of Analysis ToolPak to add a check mark/OK

If you need help doing this, ask your favorite “computer techie” for help.
You are now ready to skip ahead to Sect. [6.5.2](#)

6.5.1.3 Installing the Data Analysis ToolPak into Excel 2013

Open a new Excel spreadsheet

Click on: Data (at the top of your screen)

Look at the top of your monitor screen. Do you see the words: “Data Analysis” at the far right of the screen? If you do, the Data Analysis ToolPak for Excel 2013 was correctly installed when you installed Office 2013, and you should skip ahead to Sect. 6.5.2.

If the words: “Data Analysis” are not at the top right of your monitor screen, then the ToolPak component of Excel 2013 was not installed when you installed Office 2013 onto your computer. If this happens, you need to follow these steps:

File

Options (bottom left of screen)

Note: This creates a dialog box with “Excel Options” at the top left of the box

Add-Ins (on left of screen)

Manage: Excel Add-Ins (at the bottom of the dialog box)

Go (at bottom center of dialog box)

Highlight: Analysis ToolPak (in the Add-Ins dialog box)

Put a check mark to the left of Analysis Toolpak

OK (at the right of this dialog box)

Data

You now should have the words: “Data Analysis” at the top right of your screen to show that this feature has been installed correctly

If you get a prompt asking you for the “installation CD,” put this CD in the CD drive and click on: OK

Note: If these steps do not work, you should try these steps instead: File/Options (bottom left)/Add-ins/Analysis ToolPak/Go/click to the left of Analysis ToolPak to add a check mark/OK

If you need help doing this, ask your favorite “computer techie” for help.

You are now ready to skip ahead to Sect. 6.5.2

6.5.2 Using Excel to Find the SUMMARY OUTPUT of Regression

You have now installed *ToolPak*, and you are ready to find the regression equation for the “best-fitting straight line” through the data points by using the following steps:

Open the Excel file: *FACINGS8* (if it is not already open on your screen)

Note: If this file is already open, and there is a gray border around the chart, you need to click on any empty cell outside of the chart to deselect the chart.

Now that you have installed *Toolpak*, you are ready to find the regression equation summarizing the relationship between the number of shelf facings of Kellogg’s Corn Flakes and the sales dollars in your data set.

Remember that you gave the name: *facings* to the X data (the predictor), and the name: *sales* to the Y data (the criterion) in a previous section of this chapter (see Sect. 6.2)

Data (top of screen)

Data analysis (far right at top of screen; see Fig. 6.28)

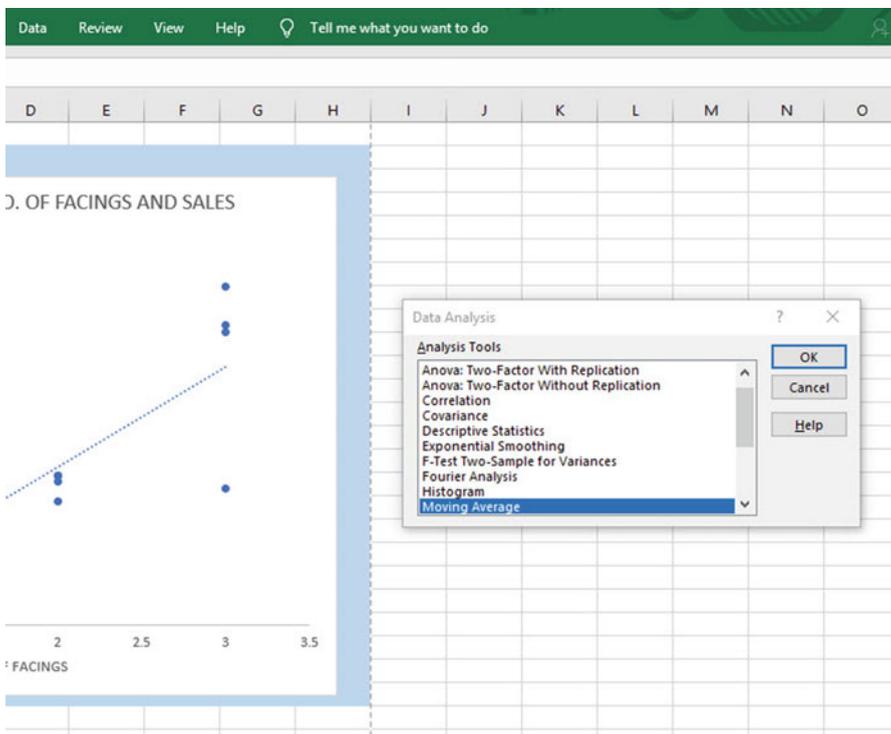
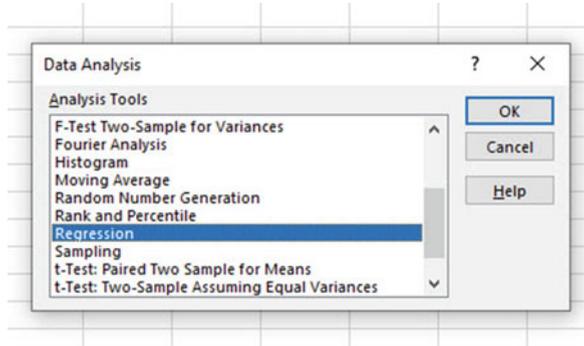


Fig. 6.28 Example of Using the Data/Data Analysis Function of Excel

Scroll down the dialog box using the down arrow and highlight: Regression (see Fig. 6.29)

Fig. 6.29 Dialogue Box for Creating the Regression Function in Excel



OK

Input Y Range: sales

Input X Range: facings

Click on the “button” to the left of Output Range to select this, and enter A44 in the box as the place on your spreadsheet to insert the Regression analysis in cell A44
OK

The *SUMMARY OUTPUT* should now be in cells: A44: I61

Widen Column A so that all of the words in the *SUMMARY OUTPUT* are readable.

Now, change the data in the following three cells to Number format (two decimal places) by first clicking on “Home” at the top left of your screen:

B47

B60

B61

Now, change the format for all other numbers that are in decimal format to number format, three decimal places.

Next, widen all columns so that all of the labels fit inside the column widths.

Then, center all numbers in their cells.

Print the file so that it fits onto one page. (*Hint: Change the scale under “Page Layout” to 70% to make it fit.*) Your file should be like the file in Fig. 6.30.



Fig. 6.30 Final Spreadsheet of Correlation and Simple Linear Regression including the SUMMARY OUTPUT for the Data

Save the resulting file as: FACINGS9

Note the following problem with the summary output.

Whoever wrote the computer program for this version of Excel made a mistake and gave the name: "Multiple R" to cell A47.

This is not correct. Instead, cell A47 should say: “correlation r ” since this is the notation that we are using for the correlation between X and Y which is $r = 0.83$.

You can now use your printout of the regression analysis to find the regression equation that is the best-fitting straight line through the data points.

But first, let’s review some basic terms.

6.5.2.1 Finding the y-Intercept, a , of the Regression Line

The point on the y -axis that the regression line would intersect the y -axis if it were extended to reach the y -axis is called the “ y -intercept” and *we will use the letter “ a ” to stand for the y -intercept of the regression line.* The y -intercept on the SUMMARY OUTPUT on the previous page is -0.65 and appears in cell B60 (note the minus sign). This means that if you were to draw an imaginary line continuing down the regression line toward the y -axis that this imaginary line would cross the y -axis at -0.65 . This is why a is called the “ y -intercept.”

6.5.2.2 Finding the Slope, b , of the Regression Line

The “tilt” of the regression line is called the “slope” of the regression line. It summarizes to what degree the regression line is either above or below a horizontal line through the data points. If the correlation between X and Y were zero, the regression line would be exactly horizontal to the X -axis and would have a zero slope.

If the correlation between X and Y is positive, the regression line would “slope upward to the right” above the X -axis. Since the regression line in Fig. 6.30 slopes upward to the right, the slope of the regression line is $+1.54$ as given in cell B61. *We will use the notation “ b ” to stand for the slope of the regression line.* (Note that Excel calls the slope of the line: “ X Variable 1” in the Excel printout.)

Since the correlation between the number of facings and the weekly sales dollars was $+0.83$, you can see that the regression line for these data “slopes upward to the right” through the data. Note that the SUMMARY OUTPUT of the regression line in Fig. 6.30 gives a correlation, r , of $+0.83$ in cell B47.

If the correlation between X and Y were negative, the regression line would “slope down to the right” above the X -axis. This would happen whenever the correlation between X and Y is a negative correlation that is between zero and minus one (0 and -1).

6.5.3 Finding the Equation for the Regression Line

To find the regression equation for the straight line that can be used to predict weekly sales from the number of facings, we only need two numbers in the SUMMARY OUTPUT in Fig. 6.30: B60 and B61.

$$\text{The format for the regression line is: } Y = a + bX \quad (6.3)$$

where $a = \text{the } y\text{-intercept}$ (-0.65 in our example in cell B60)

and $b = \text{the slope of the line}$ ($+1.54$ in our example in cell B61)

Therefore, the equation for the best-fitting regression line for our example is:

$$Y = a + bX$$

$$Y = -0.65 + 1.54X$$

Remember that Y is the weekly sales (\$000) that we are trying to predict, using the number of facings as the predictor, X .

Let's try an example using this formula to predict the weekly sales.

6.5.4 Using the Regression Line to Predict the y -Value for a Given x -Value

Objective: Find the weekly sales predicted from *one facing* of Kellogg's Corn Flakes on the supermarket shelf.

Since the number of facings is one (i.e., $X = 1$), substituting this number into our regression equation gives:

$$Y = -0.65 + 1.54(1)$$

$$Y = -0.65 + 1.54$$

$$Y = 0.89$$

Important note: If you look at your chart, if you go directly upwards from one facing until you hit the regression line, you see that you hit this line just under the number 1 on the y -axis to the left (actually, it is 0.89), the result above for predicting sales from one shelf facing.

But since weekly sales are recorded in thousands of dollars (\$000), we need to multiply our answer above by 1000 to find the weekly sales figure.

When we do that, this gives an estimated weekly sales of \$890 (0.89×1000) when we use one facing of this cereal.

Now, let's do a second example and predict what the weekly sales figure would be if we used three facings of Kellogg's Corn Flakes on the supermarket shelf.

$$Y = -0.65 + 1.54 X$$

$$Y = -0.65 + 1.54 (3)$$

$$Y = -0.65 + 4.62$$

$$Y = 3.97$$

Important note: If you look at your chart, if you go directly upwards from three facings until you hit the regression line, you see that you hit this line just under the number 4 on the y-axis to the left (actually it is 3.97), the result above for predicting sales from three shelf facings.

But since weekly sales are recorded in thousands of dollars (\$000), we need to multiply our answer above by 1000 to find the weekly sales figure.

When we do that, this gives an estimated weekly sales of \$3970 when we use three facings of the cereal.

For a more detailed discussion of regression, see Black (2010).

6.6 Adding the Regression Equation to the Chart

Objective: To Add the Regression Equation to the Chart

If you want to include the regression equation within the chart next to the regression line, you can do that, but a word of caution first.

Throughout this book, we are using the regression equation for one predictor and one criterion to be the following:

$$Y = a + b X \tag{6.3}$$

where a = y-intercept and
 b = slope of the line

See, for example, the regression equation in Sect. 6.5.3 where the y-intercept was $a = -0.65$ and the slope of the line was $b = +1.54$ to generate the following regression equation:

$$Y = -0.65 + 1.54 X$$

However, Excel 2019 uses a slightly different regression equation (which is logically identical to the one used in this book) when you add a regression equation to a chart:

$$Y = b X + a \tag{6.4}$$

where $a = y$ -intercept and $b =$ slope of the line

Note that this equation is identical to the one we are using in this book with the terms arranged in a different sequence.

For the example we used in Sect. 6.5.3, Excel 2019 would write the regression equation on the chart as:

$$Y = 1.54 X - 0.65$$

This is the format that will result when you add the regression equation to the chart using Excel 2019 using the following steps:

Open the file: FACINGS9 (that you saved in Sect. 6.5.2)

Click just *inside* the outer border of the chart in the top right corner to add the “border” around the chart in order to “select the chart” for changes you are about to make

Right-click on any of the data-points in the chart

Highlight: Add Trendline, and click on it to select this command

The “Linear button” near the top of the dialog box will already be selected (on its left)

Scroll down this dialog box, and click on: Display Equation on chart (near the bottom of the dialog box; see Fig. 6.31)

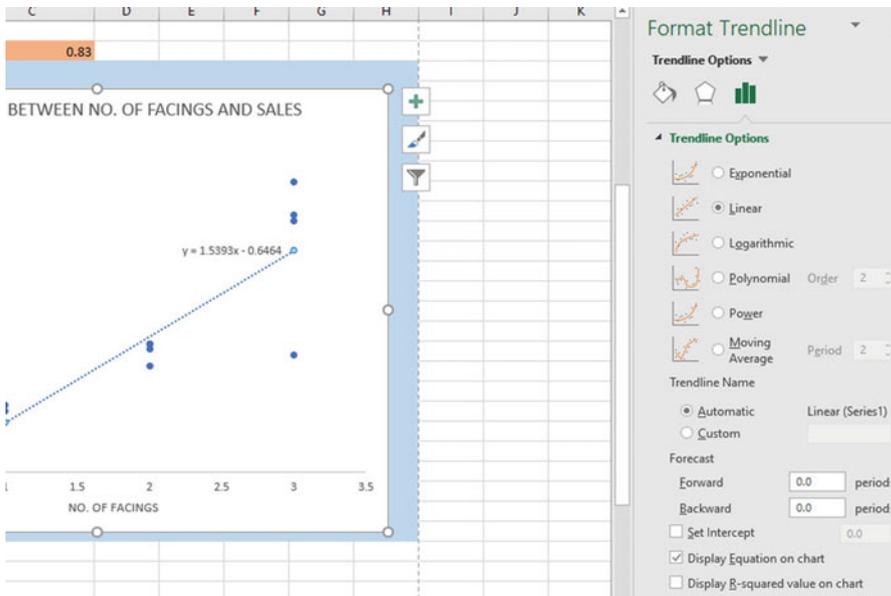


Fig. 6.31 Dialogue Box for Adding the Regression Equation to the Chart Next to the Regression Line on the Chart

Click on the X at the top right of the Format Trendline dialogue box to remove this box.

Click on any empty cell outside of the chart to deselect the chart.

Note that the regression equation on the chart is in the following form next to the regression line on the chart (see Fig. 6.32).

$$Y = 1.54X - 0.65$$

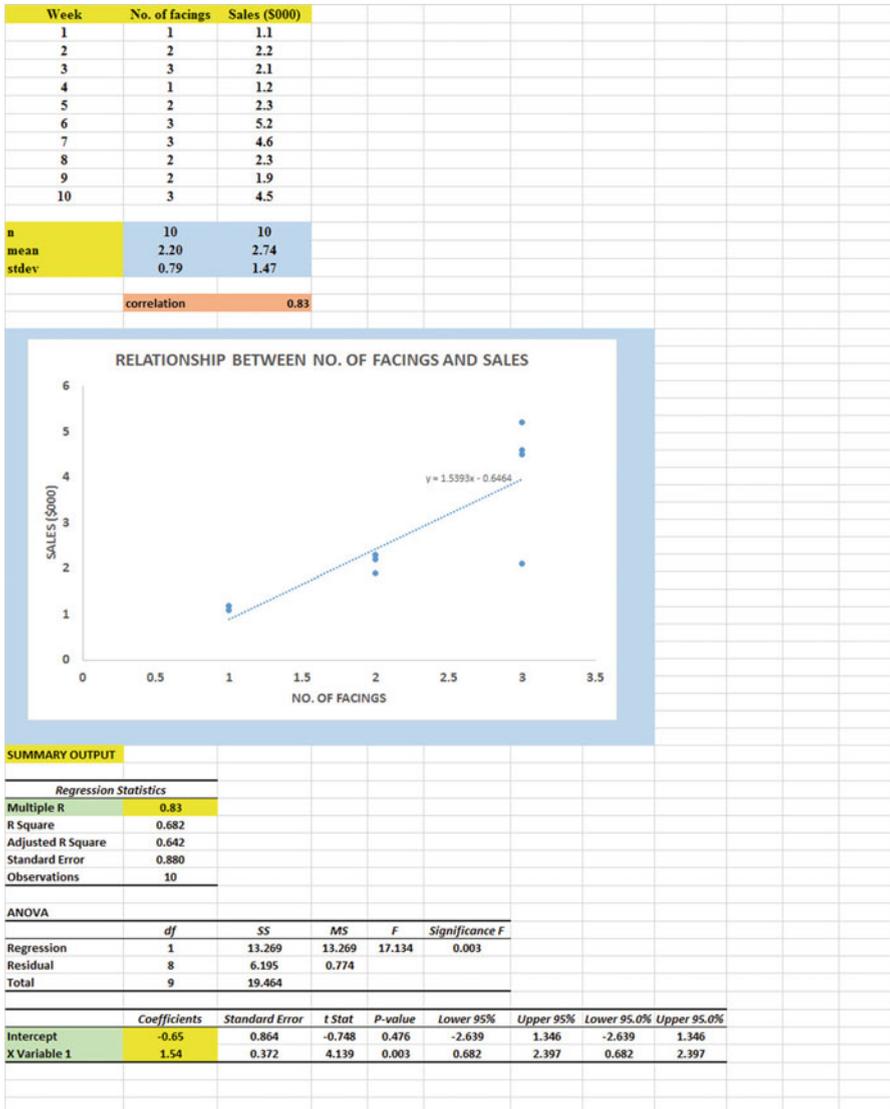


Fig. 6.32 Example of a Chart with the Regression Equation Displayed Next to the Regression Line

(Save this file as: FACINGS10, and print it out so that it fits onto one page)

6.7 How to Recognize Negative Correlations in the SUMMARY OUTPUT Table

Important note: Since Excel does not recognize negative correlations in the SUMMARY OUTPUT results, but treats all correlations as if they were positive correlations (this was a mistake made by the programmer), you need to be careful to note that there may be a negative correlation between X and Y even if the printout says that the correlation is a positive correlation.

You will know that the correlation between X and Y is a negative correlation when these two things occur:

- (1) *THE SLOPE, b, IS A NEGATIVE NUMBER. This can only occur when there is a negative correlation.*
- (2) *THE CHART CLEARLY SHOWS A DOWNWARD SLOPE IN THE REGRESSION LINE, which can only occur when the correlation between X and Y is negative.*

6.8 Printing Only Part of a Spreadsheet Instead of the Entire Spreadsheet

Objective: To print part of a spreadsheet separately instead of printing the entire spreadsheet

There will be many occasions when your spreadsheet is so large in the number of cells used for your data and charts that you only want to print part of the spreadsheet separately so that the print will not be so small that you cannot read it easily.

We will now explain how to print only part of a spreadsheet onto a separate page by using three examples of how to do that using the file, FACINGS10, that you created in Sect. 6.6: (1) printing only the table and the chart on a separate page, (2) printing only the chart on a separate page, and (3) printing only the SUMMARY OUTPUT of the regression analysis on a separate page.

Note: If the file: FACINGS10 is not open on your screen, you need to open it now.

If the “border” is around the outside of the chart, click on any white space outside of the chart to deselect the chart.

Let’s describe how to do these three goals with three separate objectives:

6.8.1 *Printing Only the Table and the Chart on a Separate Page*

Objective: To print only the table and the chart on a separate page

1. Left-click your mouse starting at the top left of the table *in cell A2* and drag the mouse *down and to the right so that all of the table and all of the chart are highlighted in light blue on your computer screen from cell A2 to cell I43* (the highlighted cells are called the “selection” cells).
2. File
 - Print
 - Print Active Sheets (hit the down arrow)
 - Print Selection
 - Print

The resulting printout should contain only the table of the data and the chart resulting from the data.

Then, click on any empty cell in your spreadsheet to deselect the table and chart.

6.8.2 *Printing Only the Chart on a Separate Page*

Objective: To print only the chart on a separate page

1. Click on any “white space” *just inside the outside border of the chart in the top right corner of the chart* to create the border around all of the borders of the chart in order to “select” the chart.
2. File
 - Print
 - Print Selected chart
 - Print selected chart (again)
 - Print

The resulting printout should contain only the chart resulting from the data.

Important note: After each time you print a chart by itself on a separate page, you should immediately click on any white space OUTSIDE the chart to remove the gray border from the border of the chart. When the gray border is on the borders of the chart, this tells Excel that you want to print only the chart by itself. Do this now!

6.8.3 *Printing Only the SUMMARY OUTPUT of the Regression Analysis on a Separate Page*

Objective: To print only the SUMMARY OUTPUT of the regression analysis on a separate page

1. Left-click your mouse at the cell just above SUMMARY OUTPUT in *cell A43* on the left of your spreadsheet and drag the mouse *down and to the right* until all of the regression output is highlighted in dark blue on your screen from A43 to I62.
2. File
Print
Print Active Sheets (hit the down arrow)
Print Selection
Print

The resulting printout should contain only the summary output of the regression analysis on a separate page.

Finally, click on any empty cell on the spreadsheet to “deselect” the regression table.

6.9 End-of-Chapter Practice Problems

1. Suppose that Schnuck’s Supermarket in Webster Groves, Missouri wanted to study the relationship between the amount of money it was spending weekly in the *Webster-Kirkwood Times* newspaper ads and the weekly sales (\$000) of its store in Webster Groves. You have decided to use a correlation and simple linear regression analysis, and to test your Excel skills, you have collected the data for this comparison. These hypothetical data appear in Fig. 6.33:

RELATIONSHIP BETWEEN NEWSPAPER AD (\$) AND WEEKLY SALES (\$000)	
NEWSPAPER AD (\$)	WEEKLY SALES (\$000)
1,480	210
2,520	170
1,540	180
2,463	220
1,810	200
1,563	205
1,465	230
2,581	170
1,712	210
1,814	240
1,653	205

Fig. 6.33 Worksheet Data for Chap. 6: Practice Problem #1

Create an Excel spreadsheet and enter the data using *NEWSPAPER AD (\$)* as the independent variable (predictor) and *WEEKLY SALES (\$000)* as the dependent variable (criterion).

Important note: When you are trying to find a correlation between two variables, it is important that you place the predictor, X, ON THE LEFT COLUMN in your Excel spreadsheet, and the criterion, Y, IMMEDIATELY TO THE RIGHT OF THE X COLUMN. You should do this every time that you want to use Excel to find a correlation between two variables to check your thinking.

- Create an Excel spreadsheet using *WEEKLY SALES (\$000)* as the criterion and *NEWSPAPER AD (\$)* as the predictor using the following format:
 - Top title: RELATIONSHIP BETWEEN NEWSPAPER AD (\$) AND WEEKLY SALES (\$000)
 - x-axis title: NEWSPAPER AD (\$)
 - y-axis title: WEEKLY SALES (\$000)
 - Re-size the chart so that it is 7 columns wide and 25 rows long
- Create the *least-squares regression line* for these data on the scatterplot.
- Use Excel's *regression* function to find the equation for the least-squares regression line for these data and display the results below the chart on your spreadsheet. Add the regression line and the regression equation to the chart.
- Use number format (two decimal places) for the correlation, the y-intercept, and the slope of the line on the SUMMARY OUTPUT, and use number format (three decimal places) for all of the other decimal figures in the SUMMARY OUTPUT.

- (e) Print the *input data and the chart* so that this information fits onto one page.
- (f) Then, print the *regression output table* so that this information fits onto a separate page.
- (g) Save the file as: NEWS3

Answer the following questions using your Excel printout:

By hand:

1. Circle and label the value of the *y-intercept* and the *slope* of the regression line on your printout.
 2. Write the *regression equation* by hand on your printout for these data
 3. Circle and label the *correlation* between the two sets of scores in the regression analysis SUMMARY OUTPUT table on your printout
 4. Underneath the regression equation you wrote by hand on your printout, use the regression equation to predict the WEEKLY SALES you would expect for a NEWSPAPER AD expense of \$2000.
 5. Read from the graph, the WEEKLY SALES you would expect for NEWSPAPER AD expense of \$1500, and write your answer on the separate page
 6. Re-save the file as: NEWS3
2. In a large engineering company, what is the relationship between the salary of engineers as a percent of the engineers' midpoint salary (position in range) and the raise given to the engineers at the last contract? The midpoint of the range of engineers' salaries is scored as 100, and each engineer's salary is than compared to that midpoint to determine what percent of that midpoint an engineer's salary represents. The resulting number is called "position in range." Engineers whose salaries are below the midpoint have a score less than 100, and engineers whose salaries are above the midpoint have a score greater than 100. Suppose that you wanted to study this question. Analyze the hypothetical data that are given in Fig. 6.34.

COMPANY XYZ	
Question:	Is there a relationship between the salary of engineers as a percent of the engineers' midpoint salary (position in range) and the raise given to the engineers at the last contract?
POSITION IN RANGE	PERCENT RAISE
83	5.5
90	5.0
100	3.0
110	1.5
86	4.0
97	3.5
102	4.0
107	1.5
112	2.0
114	2.5
116	1.5

Fig. 6.34 Worksheet Data for Chap. 6: Practice Problem #2

Create an Excel spreadsheet, and enter the data.

- (a) create an *XY scatterplot* of these two sets of data such that:
- top title: RELATIONSHIP BETWEEN POSITION IN RANGE AND PERCENT RAISE FOR ENGINEERS
 - x-axis title: POSITION IN RANGE
 - y-axis title: % RAISE
 - move the chart below the table
 - re-size the chart so that it is 7 columns wide and 25 rows long

- (b) Create the *least-squares regression line* for these data on the scatterplot.
- (c) Use Excel to run the regression statistics to find the *equation for the least-squares regression line* for these data and display the results below the chart on your spreadsheet. Add the regression equation to the chart. Use number format (two decimal places) for the correlation and number format (three decimal places) for the coefficients.

Print *just the input data and the chart* so that this information fits onto one page in portrait format.

Then, print *just the regression output table* on a separate page so that it fits onto that separate page in portrait format.

By hand:

- (d) Circle and label the value of the *y-intercept* and the *slope* of the regression line on your printout.
- (e) Write the regression equation *by hand* on your printout for these data (use three decimal places for the *y-intercept* and the *slope*).
- (f) Circle and label the *correlation* between the two sets of scores in the regression analysis summary output table on your printout.
- (g) Underneath the regression equation you wrote by hand on your printout, use the regression equation to predict the PERCENT RAISE you would expect for an engineer with a POSITION IN RANGE score of 90.
- (h) *Read from the graph*, the PERCENT RAISE you would expect for an engineer with a POSITION IN RANGE score of 110, and write your answer in the space immediately below:
-
- (i) save the file as: ENGINE3
3. Is there a relationship between the number of sales calls a sales staff make in a month on potential customers and the number of copier machines sold that month by a salesperson? Suppose that you gathered the hypothetical data given below for your sales staff for the previous month. The resulting data are presented in Fig. 6.35.

Fig. 6.35 Worksheet Data
for Chap. 6: Practice
Problem #3

No. of sales calls	No. of copiers sold
25	40
30	55
18	30
22	35
14	18
18	23
22	28
24	38
12	15
13	16
18	25
22	28
25	36

Create an Excel spreadsheet and enter the data using the number of sales calls as the independent variable (predictor) and the number of copiers sold last month by each salesperson as the dependent variable (criterion).

- Use Excel's `=correl` function to find the correlation between these two sets of scores, and round off the result to two decimal places.
- create an *XY scatterplot* of these two sets of data such that:
 - top title: RELATIONSHIP BETWEEN NO. OF SALES CALLS AND COPIERS SOLD
 - x-axis title: NO. OF SALES CALLS
 - y-axis title: NO. OF COPIERS SOLD
 - move the chart below the table
 - re-size the chart so that it is 7 columns wide and 25 rows long
- Create the *least-squares regression line* for these data on the scatterplot.
- Use Excel to run the regression statistics to find the *equation for the least-squares regression line* for these data and display the results below the chart on your spreadsheet. Use number format (two decimal places) for the correlation and for the coefficients
- Print just the input data and the chart so that this information fits onto one page. Then, print the regression output table on a separate page so that it fits onto that separate page.
- save the file as: copier4

Answer the following questions using your Excel printout:

1. What is the correlation between the number of sales calls and the number of copiers sold?
2. What is the y-intercept?
3. What is the slope of the line?
4. What is the regression equation?
5. Use the regression equation to predict the number of copiers sold you would expect for a salesperson who made 25 sales calls last month. Show your work on a separate sheet of paper.

References

- Black, K. Business Statistics: For Contemporary Decision Making (6th ed.). Hoboken, NJ: John Wiley & Sons, Inc., 2010.
- Levine, D.M., Stephan, D.F., Krehbiel, T.C., and Berenson, M.L. Statistics for Managers Using Microsoft Excel (6th ed.). Boston, MA: Prentice Hall/Pearson, 2011.
- Zikmund, W.G. and Babin, B.J. Exploring Marketing Research (10th ed.). Mason, OH: South-Western Cengage Learning, 2010.

Chapter 7

Multiple Correlation and Multiple Regression



There are many times in business when you want to predict a criterion, Y , but you want to find out if you can develop a better prediction model by using *several predictors* in combination (e.g. X_1, X_2, X_3 , etc.) instead of a single predictor, X .

The resulting statistical procedure is called “multiple correlation” because it uses two or more predictors in combination to predict Y , instead of a single predictor, X . Each predictor is “weighted” differently based on its separate correlation with Y and its correlation with the other predictors. The job of multiple correlation is to produce a regression equation that will weight each predictor differently and in such a way that the combination of predictors does a better job of predicting Y than any single predictor by itself. We will call the multiple correlation: R_{xy} .

IMPORTANT NOTE: *You will remember from Chap. 6 (see Sect. 6.1) that the correlation, r , ranges from -1 to $+1$, and, therefore, can be a negative number.*

However, the multiple correlation, R_{xy} , only ranges from zero to $+1$ (0 to $+1$), and can never be negative! It is very important that you remember this fact.

You will recall (see Sect. 6.5.3) that the regression equation that predicts Y when only one predictor, X , is used is:

$$Y = a + bX \tag{7.1}$$

7.1 Multiple Regression Equation

The multiple regression equation follows a similar format and is:

$$Y = a + b_1X_1 + b_2X_2 + b_3X_3 + \text{etc. depending on the number of predictors used} \tag{7.2}$$

The “weight” given to each predictor in the equation is represented by the letter “b” with a subscript to correspond to the same subscript on the predictors.

Important note: In order to do multiple regression, you need to have installed the “Data Analysis ToolPak” that was described in Chap. 6 (see Sect. 6.5.1). If you did not install this, you need to do so now.

Let’s try a practice problem.

Suppose that you have been hired by a car rental company to see if you could predict annual sales based on the number of cars that a rental car company has in its fleet and the number of locations where you can rent that company’s cars in the U.S.

Let’s use the following notation:

Y	Annual Sales (in millions of dollars)
X ₁	No. of cars in the fleet (in thousands of cars)
X ₂	No. of locations in the U.S.

Suppose, further, that this rental car company supplied you with the following hypothetical data summarizing its performance along with the performance of its competitors (see Fig. 7.1):

CAR RENTAL COMPANIES		
Y	X1	X2
SALES (\$millions)	NO. OF CARS (000)	NO. OF LOCATIONS
1070	120	152
1460	180	1120
1480	85	1032
552	92	440
2105	315	2587
308	71	1697
2380	221	1153
1140	142	922
43	25	105
154	35	1483
72	15	442
81	18	251
333	42	465
91	15	492
147	18	44

Fig. 7.1 Worksheet Data for Rental Car Companies (Practical Example)

Create an Excel spreadsheet for these data using the following cell reference:

A3: CAR RENTAL COMPANIES

A5: Y

A6: SALES (\$millions)

A7: 1070

B5: X1

B6: NO. OF CARS (000)

B7: 120

C5: X2

C6: NO. OF LOCATIONS

C7: 152

Next, change the column width to match the above table, and change all figures to number format (zero decimal places).

Now, fill in the additional data in the chart such that:

A21: 147

B21: 18

C21: 44 (Then, center the information in all cells of your table.)

Important note: Be sure to double-check all of your numbers in your table to be sure that they are correct, or your spreadsheets will be incorrect.

Save this file as: RENTAL5

Before we do the multiple regression analysis, we need to try to make one important point very clear:

Important: When we used one predictor, X, to predict one criterion, Y, we said that you need to make sure that the X variable is ON THE LEFT in your table, and the Y variable is ON THE RIGHT in your table so that you know which variable is the predictor, and which variable is the criterion (see Sect. 6.3).

However, in multiple regression, you need to follow this rule which is exactly the opposite:

When you use several predictors in multiple regression, it is essential that the criterion you are trying to predict, Y, be ON THE FAR LEFT, and all of the predictors are TO THE RIGHT of the criterion, Y, in your table so that you know which variable is the criterion, Y, and which variables are the predictors.

Notice in the table above, that the criterion Y (SALES) is on the far left of the table, and the two predictors (NO. OF CARS and NO. OF LOCATIONS) are to the right of the criterion variable. You must follow this rule or your regression equation will be completely wrong.

7.2 Finding the Multiple Correlation and the Multiple Regression Equation

Objective: To find the multiple correlation and multiple regression equation using Excel.

You do this by the following commands:

Data

Click on: Data Analysis (far right top of screen)

Regression (scroll down to this in the box; see Fig. 7.2)

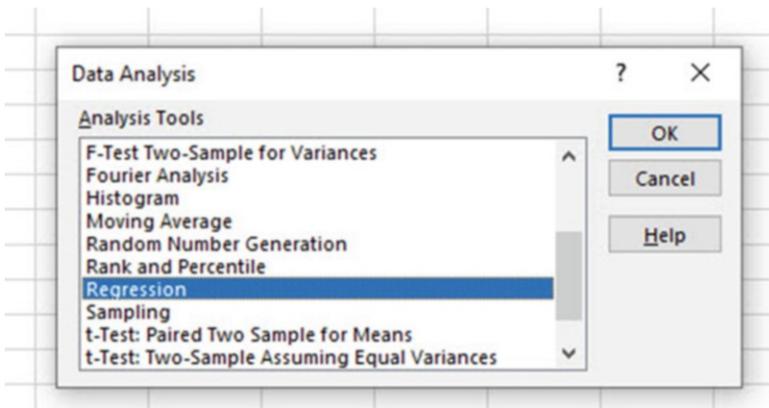


Fig. 7.2 Dialogue Box for Regression Function

OK

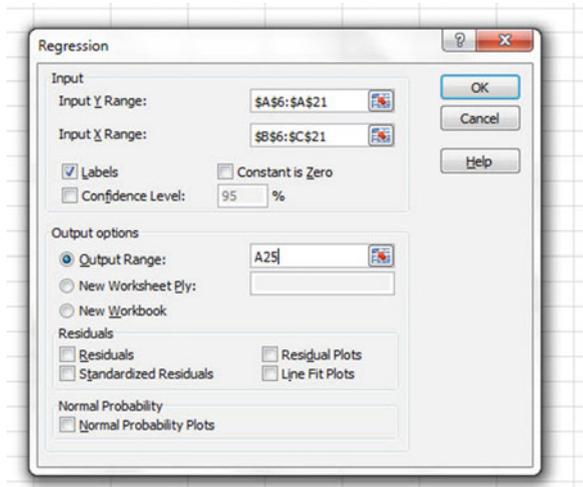
Input Y Range: A6:A21

Input X Range: B6:C21

Click on the Labels box to *add a check mark* to it (because you have included the column labels in row 6)

Output Range (click on the button to its left, and enter): A25 (see Fig. 7.3)

Fig. 7.3 Dialogue Box for Regression of Car Rental Companies Data



Important note: Excel automatically assigns a dollar sign \$ in front of each column letter and each row number so that you can keep these ranges of data constant for the regression analysis.

OK (see Fig. 7.4 to see the resulting SUMMARY OUTPUT)

	A	B	C	D	E	F	G	H	I	J
21	147	18	44							
22										
23										
24										
25	SUMMARY OUTPUT									
26										
27	<i>Regression Statistics</i>									
28	Multiple R		0.93							
29	R Square		0.86							
30	Adjusted R Square		0.83							
31	Standard Error		321.49							
32	Observations		15							
33										
34	ANOVA									
35		<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>				
36	Regression	2	7510945.33	3755473	36.33	8.10477E-06				
37	Residual	12	1240299.61	103358						
38	Total	14	8751244.93							
39										
40		<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>	<i>Lower 95.0%</i>	<i>Upper 95.0%</i>	
41	Intercept	53.55	133.20	0.40	0.69	-236.66	343.76	-236.66	343.76	
42	NO. OF CARS (000)	9.09	1.34	6.78	0.00	6.17	12.01	6.17	12.01	
43	NO. OF LOCATIONS	-0.17	0.17	-0.98	0.34	-0.53	0.20	-0.53	0.20	
44										

Fig. 7.4 Regression SUMMARY OUTPUT of Car Rental Companies Data

Next, format the following four cells in Number format (two decimal places):

- B28
- B41
- B42
- B43

Note that both the input Y Range and the Input X Range above both include the label at the top of the columns. Format all other decimal figures to two decimals.

Re-save the file as: RENTAL5

Now, print the file so that it fits onto one page by changing the scale to 60% size. The resulting regression analysis is given in Fig. 7.5.

CAR RENTAL COMPANIES		
Y	X1	X2
SALES (\$millions)	NO. OF CARS (000)	NO. OF LOCATIONS
1070	120	152
1460	180	1120
1480	85	1032
552	92	440
2105	315	2587
308	71	1697
2380	221	1153
1140	142	922
43	25	105
154	35	1483
72	15	442
81	18	251
333	42	465
91	15	492
147	18	44

SUMMARY OUTPUT	
Regression Statistics	
Multiple R	0.93
R Square	0.86
Adjusted R Square	0.83
Standard Error	321.49
Observations	15

ANOVA					
	df	SS	MS	F	Significance F
Regression	2	7510945.33	3755473	36.33	8.10477E-06
Residual	12	1240299.61	103358		
Total	14	8751244.93			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	53.55	133.20	0.40	0.69	-236.66	343.76	-236.66	343.76
NO. OF CARS (000)	9.09	1.34	6.78	0.00	6.17	12.01	6.17	12.01
NO. OF LOCATIONS	-0.17	0.17	-0.98	0.34	-0.53	0.20	-0.53	0.20

Fig. 7.5 Final Spreadsheet for Car Rental Companies Regression Analysis

Once you have the SUMMARY OUTPUT, you can determine the multiple correlation and the regression equation that is the best-fit line through the data points using NO. OF CARS (000) and NO. OF LOCATIONS as the two predictors, and SALES (\$millions) as the criterion.

Note on the SUMMARY OUTPUT where it says: “Multiple R” in cell A28. This term is correct since this is the term Excel uses for the multiple correlation, which is +0.93. This means, that from these data, that the combination of NO. OF CARS and NO. OF LOCATIONS together form a very strong positive relationship in predicting Annual Sales.

To find the regression equation, *notice the coefficients at the bottom of the SUMMARY OUTPUT:*

<i>Intercept: a (this is the y-intercept)</i>	53.55
<i>NO. OF CARS (000): b1</i>	9.09
<i>NO. OF LOCATIONS: b2</i>	-0.17

Since the general form of the multiple regression equation is:

$$Y = a + b_1X_1 + b_2X_2 \quad (7.2)$$

we can now write the multiple regression equation for these data:

$$Y = 53.55 + 9.09X_1 - 0.17X_2$$

7.3 Using the Regression Equation to Predict Annual Sales

Objective: To find the predicted annual sales for a rental car company that has 80,000 cars and 900 locations.

Note that X_1 (NO. OF CARS) is measured in thousands of cars in the original data set. This means, that for our example, that 80,000 cars would become just 80, since 80 is 80,000 measured in thousands of cars. Plugging these two numbers into our regression equation gives us:

$$Y = 53.55 + 9.09 (80) - 0.17 (900)$$

$$Y = 53.55 + 727.2 - 153$$

$$Y = 627.75$$

But, since Annual Sales are measured in millions of dollars in the original data set, we have to convert this figure to millions of dollars. Therefore, the predicted annual sales for a rental car company that has 80,000 cars and 900 locations where customers can rent their cars is:

\$ 627,750,000 or \$ 627.75 million

If you want to learn more about the theory behind multiple regression, see Keller (2009).

7.4 Using Excel to Create a Correlation Matrix in Multiple Regression

The final step in multiple regression is to find the correlation between all of the variables that appear in the regression equation.

In our example, this means that we need to find the correlation between each of the three pairs of variables:

- (1) number of cars and sales
- (2) number of locations and sales
- (3) number of cars and number of locations

To do this, we need to use Excel to create a “correlation matrix.” This matrix summarizes the three correlations above.

Objective: To use Excel to create a correlation matrix between the three variables in this example.

To use Excel to do this, use these steps:

Data (top of screen under “Home” at the top left of screen)

Data Analysis

Correlation (scroll *up* to highlight this formula; see Fig. 7.6)

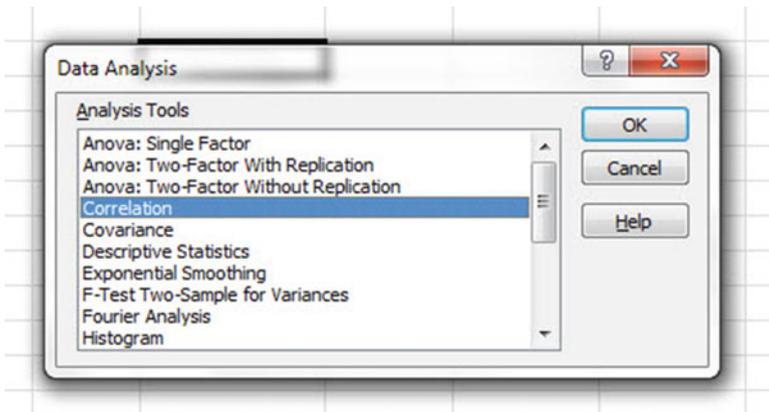


Fig. 7.6 Dialogue Box for Correlation Matrix for Car Rental Companies

OK

Input range: A6:C21

(Note that this input range includes the labels at the top of the three variables (SALES, NO. OF CARS, and NO. OF LOCATIONS) as well as all of the figures in the original data set.)

Grouped by: Columns

Put a check in the box for: Labels in the First Row (since you included the labels at the top of the columns in your input range of data above)

Output range (click on the button to its left, and enter): A47 (see Fig. 7.7)

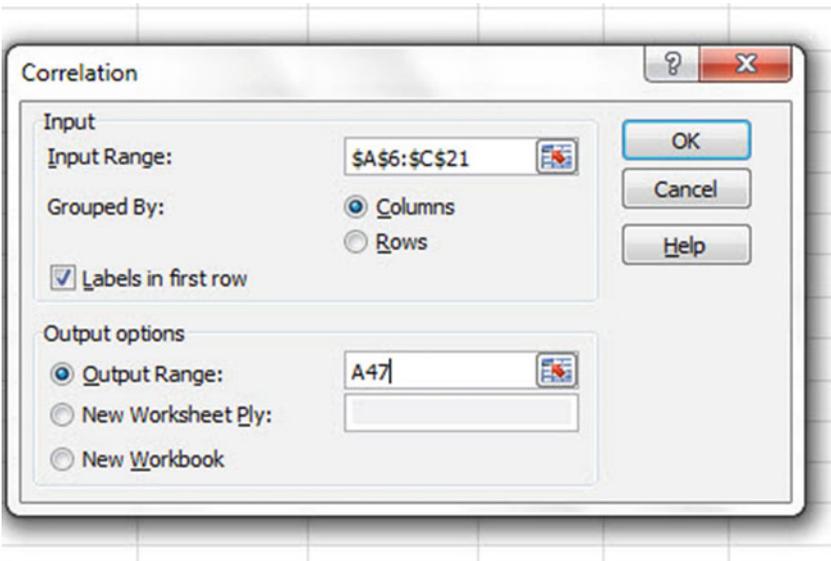


Fig. 7.7 Dialogue Box for Input/Output Range for Correlation Matrix

OK

The resulting correlation matrix appears in A47:D50 (See Fig. 7.8).

	SALES (\$millions)	NO. OF CARS (000)	NO. OF LOCATIONS
SALES (\$millions)	1		
NO. OF CARS (000)	0.920235314	1	
NO. OF LOCATIONS	0.562140716	0.694488326	1

Fig. 7.8 Resulting Correlation Matrix for Rental Car Companies Data

Next, format the three numbers in the correlation matrix that are in decimals to two decimal places. And, also, make column D wider so that the Number of Locations label fits inside cell D47. Center all numbers in the correlation matrix.

Save this Excel file as: RENTAL6

The final spreadsheet for these Car Rental Companies appears in Fig. 7.9.

CAR RENTAL COMPANIES		
Y	X1	X2
SALES (\$millions)	NO. OF CARS (000)	NO. OF LOCATIONS
1070	120	152
1460	180	1120
1480	85	1032
552	92	440
2105	315	2587
308	71	1697
2380	221	1153
1140	142	922
43	25	105
154	35	1483
72	15	442
81	18	251
333	42	465
91	15	492
147	18	44

SUMMARY OUTPUT	
Regression Statistics	
Multiple R	0.93
R Square	0.86
Adjusted R Square	0.83
Standard Error	321.49
Observations	15

ANOVA					
	df	SS	MS	F	Significance F
Regression	2	7510945.33	3755472.663	36.33	8.10477E-06
Residual	12	1240299.61	103358.3006		
Total	14	8751244.93			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	53.55	133.20	0.40	0.69	-236.66	343.76	-236.66	343.76
NO. OF CARS (000)	9.09	1.34	6.78	0.00	6.17	12.01	6.17	12.01
NO. OF LOCATIONS	-0.17	0.17	-0.98	0.34	-0.53	0.20	-0.53	0.20

	SALES (\$millions)	NO. OF CARS (000)	NO. OF LOCATIONS
SALES (\$millions)	1		
NO. OF CARS (000)	0.92	1	
NO. OF LOCATIONS	0.56	0.69	1

Fig. 7.9 Final Spreadsheet for Car Rental Companies Regression and the Correlation Matrix

Note that the number “1” along the diagonal of the correlation matrix means that the correlation of each variable with itself is a perfect, positive correlation of 1.0. *Correlation coefficients are always expressed in just two decimal places.*

You are now ready so read the correlation between the three pairs of variables:

- The correlation between NO. OF CARS and SALES is:* +.92
- The correlation between NO. OF LOCATIONS and SALES is:* +.56
- The correlation between NO. OF CARS and NO. OF LOCATIONS is:* +.69

This means that the better predictor of sales is NO. OF CARS with a correlation of +.92. Adding the second predictor variable, NO. OF LOCATIONS, improved

the prediction by only .01 to 0.93, and was, therefore, not worth the extra effort. NO. OF CARS is an excellent prediction of ANNUAL SALES all by itself.

If you want to learn more about the correlation matrix, see Levine et al. (2011).

7.5 End-of-Chapter Practice Problems

1. The Graduate Record Examinations (GRE) are frequently used to predict the first-year GPA of students in an MBA program.

The Graduate Record Examinations (GRE) are a standardized test that is an admissions requirement for many U.S. graduate schools that offer an MBA degree. The GRE is intended to measure general academic preparedness, regardless of specialization field. The GRE test produces three subtest scores: (1) GRE VERBAL REASONING (scale 130–170), (2) GRE QUANTITATIVE REASONING (scale 130–170), and (3) ANALYTICAL WRITING (scale 0–6).

Suppose that you have been asked by a director of an MBA program to find out the relationship between these variables based on last year’s entering graduate class and the ability of the GRE to predict first-year grade-point average (GPA).

You have decided to use the three subtest scores as the predictors, X_1 , X_2 , and X_3 and the first-year grade-point average (FIRST-YEAR GPA) as the criterion, Y . To test your Excel skills, you have randomly selected a small group of students from last year’s entering MBA class, and have recorded their scores on these variables.

But, suppose, that you want to find out what would happen if you added undergraduate GPA as a fourth predictor. What would be the multiple correlation?

Let’s find out what happens when you use the hypothetical data that is presented in Fig. 7.10 that includes undergraduate GPA as a fourth predictor of first-year GPA for students in an MBA program.

GRADUATE RECORD EXAMINATIONS (GRE)				
How well does the GRE predict first-year GPA in an MBA program?				
FIRST-YEAR GPA	GRE VERBAL	GRE QUANTITATIVE	GRE WRITING	UNDERGRAD GPA
3.25	160	161	5	3.40
3.42	156	158	4	3.15
2.85	156	157	2	3.05
2.65	154	153	1	2.55
3.65	166	166	6	3.25
3.16	159	160	3	3.20
3.56	166	163	4	3.66
2.35	155	154	2	2.55
2.86	153	154	3	2.85
2.95	158	157	4	2.80
3.15	158	159	4	3.05
3.45	160	160	5	3.44

Fig. 7.10 Worksheet Data for Chap. 7 Practice Problem #1

- (a) Create an Excel spreadsheet using FIRST-YEAR GPA as the criterion (Y), and the other variables as the four predictors of this criterion.
- (b) Use Excel's *multiple regression* function to find the relationship between these variables and place it below the table.
- (c) Use number format (two decimal places) for the multiple correlation on the Summary Output, use number format (three decimal places) for the coefficients, and four decimal places for all other decimal figures in the SUMMARY OUTPUT.
- (d) Print the table and regression results below the table so that they fit onto one page.
- (e) By hand on this printout, *circle and label*:
 - (1a) multiple correlation R_{xy}
 - (2b) coefficients for the y-intercept, GRE VERBAL, GRE QUANTITATIVE, GRE WRITING, AND UNDERGRAD GPA
- (f) Save this file as: GRE24
- (g) Now, go back to your Excel file and create a correlation matrix for these five variables, and place it underneath the SUMMARY OUTPUT. *Change each correlation to just two decimals.* Save this file again as: GRE24
- (h) Now, print out *just this correlation matrix in portrait mode* on a separate sheet of paper.

Answer the following questions using your Excel printout:

1. What is the multiple correlation R_{xy} ?
2. What is the y-intercept a ?
3. What is the coefficient for GRE VERBAL b_1 ?
4. What is the coefficient for GRE QUANTITATIVE b_2 ?
5. What is the coefficient for GRE WRITING b_3 ?
6. What is the coefficient for UNDERGRAD GPA b_4 ?
7. What is the multiple regression equation?
8. Underneath this regression equation by hand, predict the FIRST-YEAR GPA you would expect for a GRE VERBAL score of 159, a GRE QUANTITATIVE score of 154, A GRE WRITING score of 4, and an UNDERGRAD GPA of 3.05.

Answer the following questions using your Excel printout. Be sure to include the plus or minus sign for each correlation:

9. What is the correlation between UNDERGRAD GPA and FIRST-YEAR GPA?
10. What is the correlation between UNDERGRAD GPA and GRE VERBAL?
11. What is the correlation between UNDERGRAD GPA and GRE QUANTITATIVE?
12. What is the correlation between UNDERGRAD GPA and GRE WRITING?

13. Discuss which of the four predictors is the best predictor of FIRST-YEAR GPA.
 14. Explain in words how much better the four predictor variables combined predict FIRST-YEAR GPA than the best single predictor by itself.
2. The Graduate Management Admission Test (GMAT) is a three-and-a-half hour exam that is accepted by almost 6000 Business and Management programs in more than 80 countries as part of the admission application for people who want to obtain a graduate degree. This test is taken by more than 200,000 applicants each year. Suppose that a major university that offers an M.A. in Human Resources Management requires a GMAT score as part of the application process to this program, wants to know how well GMAT scores of applicants predict their Grade-Point Average (GPA) at the end of their first year of graduate school. The GMAT has four subtest scores: (1) Verbal (score range 0–60), (2) Quantitative (score range 0–60), (3) Analytical writing (score range 0–6 in 0.5 intervals), and (4) Integrated Reasoning (score range 1–8) You have decided to use these four subtest scores as predictors of first-year GPA, and to check your skills in Excel, you have created the hypothetical data given in Fig. 7.11.

GRADUATE MANAGEMENT ADMISSION TEST (GMAT)				
How well does the GMAT predict first-year GPA in an HRM program?				
FIRST-YEAR GPA	VERBAL	QUANTITATIVE	ANALYTICAL WRITING	INTEGRATED REASONING
3.25	50	45	4.0	4
3.67	48	48	4.5	6
2.80	35	51	5.0	5
3.05	41	50	5.5	4
3.45	51	49	4.0	3
3.33	48	45	3.0	7
2.75	46	51	4.5	8
2.95	45	48	5.5	5
2.60	40	51	6.0	6
3.67	50	50	4.5	4
3.75	46	48	3.0	7
3.42	46	46	4.0	6
3.15	42	48	5.0	7
3.26	38	49	4.0	5
2.96	41	51	5.5	4

Fig. 7.11 Worksheet Data for Chap. 7: Practice Problem #2

- (a) create an Excel spreadsheet using FIRST-YEAR GPA as the criterion (Y), and the other variables as the four predictors of this criterion ($X_1 = \text{VERBAL}$, $X_2 = \text{QUANTITATIVE}$, $X_3 = \text{ANALYTICAL WRITING}$, and $X_4 = \text{INTEGRATED REASONING}$).
- (b) Use Excel's *multiple regression* function to find the relationship between these five variables and place the SUMMARY OUTPUT below the table.

- (c) Use number format (two decimal places) for the multiple correlation on the Summary Output, and use three decimal places for the coefficients in the SUMMARY OUTPUT.
- (d) Save the file as: GMAT7
- (e) Print the table and regression results below the table so that they fit onto one page.

Answer the following questions using your Excel printout:

1. What is the multiple correlation R_{xy} ?
 2. What is the y-intercept a ?
 3. What is the coefficient for VERBAL, b_1 ?
 4. What is the coefficient for QUANTITATIVE, b_2 ?
 5. What is the coefficient for ANALYTICAL WRITING, b_3 ?
 6. What is the coefficient for INTEGRATED REASONING, b_4 ?
 7. What is the multiple regression equation?
 8. Predict the FIRST-YEAR GPA you would expect for a VERBAL score of 48, a QUANTITATIVE SCORE OF 46, an ANALYTICAL WRITING SCORE of 4.5, and an INTEGRATED REASONING SCORE OF 6.
- (f) Now, go back to your Excel file and create a correlation matrix for these five variables, and place it underneath the SUMMARY OUTPUT.
 - (g) Re-save this file as: GMAT8
 - (h) Now, print out *just this correlation matrix* on a separate sheet of paper.

Answer to the following questions using your Excel printout. (Be sure to include the plus or minus sign for each correlation):

9. What is the correlation between VERBAL and FIRST-YEAR GPA?
 10. What is the correlation between QUANTITATIVE and FIRST-YEAR GPA?
 11. What is the correlation between ANALYTICAL WRITING and FIRST-YEAR GPA?
 12. What is the correlation between INTEGRATED REASONING and FIRST-YEAR GPA?
 13. What is the correlation between VERBAL and QUANTITATIVE?
 14. What is the correlation between QUANTITATIVE and ANALYTICAL WRITING?
 15. What is the correlation between ANALYTICAL WRITING and INTEGRATED REASONING?
 16. What is the correlation between QUANTITATIVE and INTEGRATED REASONING?
 17. Discuss which of the four predictors is the best predictor of FIRST-YEAR GPA.
 18. Explain in words how much better the four predictor variables combined predict FIRST-YEAR GPA than the best single predictor by itself.
3. Suppose that you are the Advertising Manager of a large Midwestern hardware store chain, and that you have been asked to determine the effectiveness of three

different types of your advertising dollars on the dollar sales of your stores (\$000) based on three predictors: (1) Direct Mail flyers to households that are located within three miles of the store (\$000), (2) local Billboards (\$000), and (3) local TV Ads (\$000). Suppose that you have selected a random sample of hardware stores and recorded the hypothetical data given in Fig. 7.12.

HARDWARE STORES RESULTS			
Sales (\$000)	Direct Mail (\$000)	Billboards (\$000)	TV Ads (\$000)
15.98	1.05	1.65	1.55
21.59	1.29	1.55	1.97
25.47	0.96	1.86	1.64
16.02	1.51	1.66	1.73
16.58	1.3	1.45	1.61
20.1	1.06	1.55	1.82
17.49	0.86	1.42	1.97
22.53	0.78	1.55	1.61
23.98	1.2	1.77	1.13
18.96	0.83	1.17	1.77
24.64	1.01	1.76	2.04
19.52	1.34	1.62	1.63
19.97	0.55	1.95	1.93
18.86	1.06	1.55	1.63
22.79	0.9	1.32	2.27
20.38	1.88	1.15	1.45
16.51	0.51	1.19	1.93
18.32	0.53	1.65	1.44
16.04	1.21	1.45	1.46
18.91	0.9	1.45	1.66
20.09	0.97	1.56	1.59
19.83	1.01	1.45	1.42
18.42	1.29	1.47	2.05

Fig. 7.12 Worksheet Data for Chap. 7: Practice Problem #3

- (a) create an Excel spreadsheet using Sales (\$000) as the criterion, and the other variables as the three predictors of this criterion.
- (b) Use Excel's **multiple regression** function to find the relationship between these variables and place it below the table.
- (c) Use number format (two decimal places for the multiple correlation on the Summary Output, and use two decimal places for the coefficients in the Summary output).

- (d) Print the table and regression results below the table **in portrait format** so that they fit onto one page.
- (e) By hand on this printout, **circle and label:**
- (1a) multiple correlation
 - (2b) coefficients for the y-intercept, direct mail flyers, billboard ads, and television ads.
- (f) On a separate sheet of paper by hand, write the multiple regression equation
- (g) Underneath this regression equation by hand, predict the sales you would expect for direct mail costs of \$1880, billboard ad costs of \$1150, and TV ad costs of \$1630.
- (h) Save this file as: **HARDWARE4A**
- (i) Now, go back to your Excel file and create a correlation matrix for these four variables, and place it underneath the SUMMARY OUTPUT table. Use two decimal places for all correlations.
- (j) Now, print out **just this correlation matrix in portrait mode** on a separate sheet of paper.
- (k) By hand underneath the correlation matrix printed out in part j, write the answer to the following questions (label your answer as 1a, 2b, 3c, etc.). Be sure to include the plus or minus sign for each correlation:
- (1a) What is the correlation between Direct Mail and Dollar Sales?
 - (2b) What is the correlation between Billboard Ads and Dollar Sales?
 - (3c) What is the correlation between TV ads and Dollar Sales?
 - (4d) What is the correlation between TV ads and Direct Mail?
 - (5e) What is the correlation between Billboard Ads and TV ads?
 - (6f) Discuss which predictor is the best single predictor of sales.
 - (7g) Explain, in words, how much better the three variables together predict sales than the best single predictor among the predictor variables.
 - (8h) Save the file as: **HARDWARE11**

References

- Keller, G. Statistics for Management and Economics (8th ed.). Mason, OH: South-Western Cengage Learning, 2009.
- Levine, D.M., Stephan, D.F., Krehbiel, T.C., and Berenson, M.L. Statistics for Managers using Microsoft Excel (6th ed.). Boston, MA: Prentice Hall/Pearson, 2011.

Chapter 8

One-Way Analysis of Variance (ANOVA)



So far in this 2019 Excel Guide, you have learned how to use a one-group t-test to compare the sample mean to the population mean, and a two-group t-test to test for the difference between two sample means. *But what should you do when you have more than two groups and you want to determine if there is a significant difference between the means of these groups?*

The answer to this question is: *Analysis of Variance (ANOVA)*.

The ANOVA test allows you to test for the difference between the means when you have *three or more groups* in your research study.

Important note: In order to do One-way Analysis of Variance, you need to have installed the “Data Analysis Toolpak” that was described in Chap. 6 (see Sect. 6.5.1). If you did not install this, you need to do that now.

Let’s suppose that you are interested in comparing prices between three major supermarket chains in St. Louis: (1) Dierberg’s, (2) Schnuck’s, and (3) Shop ‘n Save. Suppose, further, that you have selected the 28 specific items listed in the table below as your “market basket of products” to compare prices at these three supermarkets. You have also specified the package size of each of these items in your checklist. Item #14, for example, might be: Tide Liquid laundry detergent, 16 ounces.

Suppose that you have selected zip code 63119 in St. Louis, as this zip code has one store of each of these three supermarket chains. You drive to each of these three supermarkets in this zip code area, and you have obtained the hypothetical data given in Fig. 8.1 summarizing the prices of the items in your market basket of products:

SUPERMARKET PRICE COMPARISONS			
ITEM	DIERBERG'S	SCHNUCK'S	SHOP 'n SAVE
1	1.85	1.45	1.25
2	3.95	3.35	3.04
3	2.25	1.75	1.45
4	2.85	2.35	2.25
5	1.65	1.10	0.85
6	3.65	2.95	2.45
7	2.45	1.85	1.45
8	1.95	1.56	1.44
9	1.83	1.25	1.15
10	2.64	2.14	2.04
11	2.84	2.25	2.15
12	1.84	1.20	0.55
13	1.65	1.25	1.15
14	2.75	2.10	2.04
15	2.71	1.86	1.75
16	1.55	0.94	0.85
17	1.85	1.30	1.01
18	0.95	0.55	0.45
19	1.55	1.28	1.06
20	1.44	0.85	0.74
21	1.65	1.25	1.15
22	1.64	1.28	1.04
23	4.21	3.75	3.36
24	1.20	0.71	0.61
25	4.55	3.90	3.25
26	3.45	2.84	2.65
27	5.85	5.30	5.14
28	1.65	1.25	1.04

Fig. 8.1 Worksheet Data for Supermarket Price Comparisons (Practical Example)

Create an Excel spreadsheet for these data in this way:

B1: SUPERMARKET PRICE COMPARISON

A3: ITEM

B3: DIERBERG'S

C3: SCHNUCK'S

D3: SHOP 'n SAVE

A4: 1

B4: 1.85

Enter the other information into your spreadsheet table. When you have finished entering these data, the last cell on the left should have 28 in cell A31, and the last cell on the right should have 1.04 in cell D31. Center the numbers and labels in each of the columns. Use number format (two decimals) for all numbers.

Widen the labels and columns so that the information looks like Fig. 8.1

Important note: Be sure to double-check all of your figures in the table to make sure that they are exactly correct or you will not be able to obtain the correct answer for this problem!

Save this file as: SUPERMARKET5

8.1 Using Excel to Perform a One-Way Analysis of Variance (ANOVA)

Objective: To use Excel to perform a one-way ANOVA test.

You are now ready to perform an ANOVA test on these data using the following steps:

Data (at top of screen)

Data Analysis (far right at top of screen)

Anova: Single Factor (*scroll up to this formula and highlight it; see Fig. 8.2*)

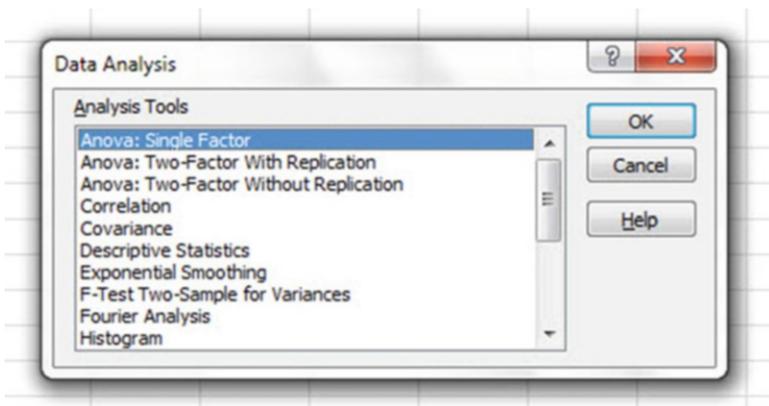


Fig. 8.2 Dialog Box for Data Analysis: Anova Single Factor

OK

Input range: B3: D31 (note that you have included in this range the column titles that are in row 3)

Important note: When you define the Input Range of the data, be sure that it includes only the data you are measuring (e.g. “Prices” in this example). Never include anything else (e.g. the ITEM numbers 1–28 in column A).

Important note: Whenever the data set has a different sample size in the groups being compared, the INPUT RANGE that you define must start at the column title of the first group on the left and go to the last column on the right and go down to the lowest row that has a figure in it in the entire data matrix so that the INPUT RANGE has the “shape” of a rectangle when you highlight it.

Grouped by: Columns

Put a check mark in: Labels in First Row

Output range (click on the button to its left): A36 (see Fig. 8.3)

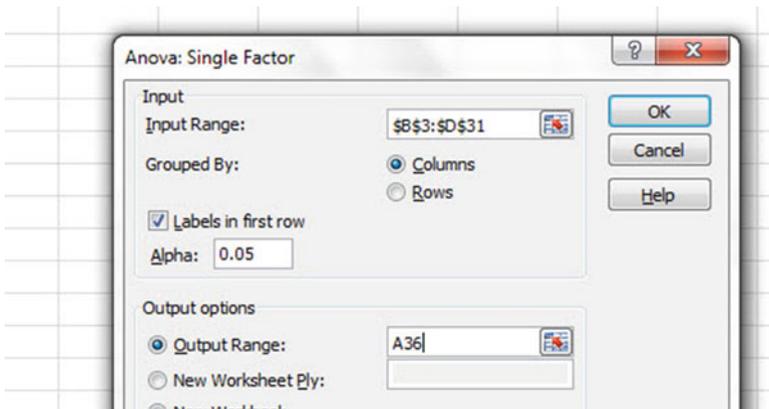


Fig. 8.3 Dialog Box for Anova: Single Factor Input/Output Range

OK

Save this file as: SUPER6

You should have generated the table given in Fig. 8.4. If you round off all figures that are in decimal format to two decimal places and center all numbers in their cells, this will make your table much easier to read.

	A	B	C	D	E	F	G	H
35								
36	Anova: Single Factor							
37								
38	SUMMARY							
39	<i>Groups</i>	<i>Count</i>	<i>Sum</i>	<i>Average</i>	<i>Variance</i>			
40	DIERBERG'S	28	68.40	2.44	1.32			
41	SCHNUCK'S	28	53.61	1.91	1.22			
42	SHOP 'n SAVE	28	47.36	1.69	1.13			
43								
44								
45	ANOVA							
46	<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F crit</i>	
47	Between Groups	8.34	2	4.17	3.40	0.04	3.11	
48	Within Groups	99.23	81	1.23				
49								
50	Total	107.57	83					

Fig. 8.4 ANOVA Results for Supermarket Price Comparisons

Print out both the data table and the ANOVA summary table so that all of this information fits onto one page. (Hint: Set the Page Layout/Fit to Scale to 85% size).

As a check on your analysis, you should have the following in these cells:

A36: Anova: Single Factor

D40: 2.44

D47: 4.17

E47: 3.40

G47: 3.11

Re-save this file as: SUPER6.

Now, let's discuss how you should interpret this table:

8.2 How to Interpret the ANOVA Table Correctly

Objective: To interpret the ANOVA table correctly

ANOVA allows you to test for the differences between means when you have three or more groups of data. This ANOVA test is called the F-test statistic, and is typically identified with the letter: F.

The formula for the F-test is this:

F = Mean Square between groups (MS_b) divided by Mean Square within groups (MS_w)

$$F = MS_b / MS_w \tag{8.1}$$

The derivation and explanation of this formula is beyond the scope of this *Excel Guide*. In this *Excel Guide*, we are attempting to teach you *how to use Excel*, and we are not attempting to teach you the statistical theory that is behind the ANOVA formulas. For a detailed explanation of ANOVA, see Weiers (2011).

Note that cell D47 contains $MS_b = 4.17$, while cell D48 contains $MS_w = 1.23$.

When you divide these two figures using their cell references in Excel, you get the answer for the F-test of 3.40 which is in cell E47. Let's discuss now the meaning of the figure: $F = 3.40$.

In order to determine whether this figure for F of 3.40 indicates a significant difference between the means of the three groups, the first step is to write the null hypothesis and the research hypothesis for the three groups of prices.

In our supermarket price comparisons, the null hypothesis states that the population means of the three groups are equal, while the research hypothesis states that the population means of the three groups are not equal and that there is, therefore, a significant difference between the population means of the three groups. Which of these two hypotheses should you accept based on the ANOVA results?

8.3 Using the Decision Rule for the ANOVA F-Test

To state the hypotheses, let's call Dierberg's as Group 1, Schnuck's as Group 2, and Shop 'n Save as Group 3. The hypotheses would then be:

$$H_0 : \mu_1 = \mu_2 = \mu_3$$

$$H_1 : \mu_1 \neq \mu_2 \neq \mu_3$$

The answer to this question is analogous to the decision rule used in this book for both the one-group t-test and the two-group t-test. You will recall that this rule (See Sects. 4.1.6 and 5.1.8) was:

If the absolute value of t is less than the critical t, you accept the null hypothesis.

or

If the absolute value of t is greater than the critical t, you reject the null hypothesis, and accept the research hypothesis.

Now, here is the decision rule for ANOVA:

Objective: To learn the decision rule for the ANOVA F-test

The decision rule for the ANOVA F-test is the following:

If the value for F is less than the critical F-value, accept the null hypothesis.

or

If the value of F is greater than the critical F-value, reject the null hypothesis, and accept the research hypothesis.

Note that Excel tells you the critical F-value in cell G47: 3.11

Therefore, our decision rule for the supermarket ANOVA test is this:

Since the value of F of 3.40 is greater than the critical F-value of 3.11, we reject the null hypothesis and accept the research hypothesis.

Therefore, our conclusion, in plain English, is:

There is a significant difference between the population means of the three supermarkets' prices.

Note that it is not necessary to take the absolute value of F of 3.40. The F-value can never be less than one, and so it can never be a negative value which requires us to take its absolute value in order to treat it as a positive value.

It is important to note that ANOVA tells us that there was a significant difference between the population means of the three groups, *but it does not tell us which pairs of groups were significantly different from each other.*

8.4 Testing the Difference Between Two Groups Using the ANOVA t-Test

To answer that question, we need to do a different test called the ANOVA t-test.

Objective: To test the difference between the means of two groups using an ANOVA t-test when the ANOVA results indicate a significant difference between the population means.

Since we have three groups of data (one group for each of the three supermarkets), we would have to perform three separate ANOVA t-tests to determine which pairs of groups were significantly different. This means that we would have to perform a separate ANOVA t-test for the following pairs of groups:

- (1) Dierberg's vs. Schnuck's
- (2) Dierberg's vs. Shop 'n Save
- (3) Schnuck's vs. Shop 'n Save

We will do just one of these pairs of tests, Dierberg's vs. Shop 'n Save, to illustrate the way to perform an ANOVA t-test comparing these two supermarkets. The ANOVA t-test for the other two pairs of groups would be done in the same way.

8.4.1 Comparing Dierberg's vs. Shop 'n Save in Their Prices Using the ANOVA t-Test

Objective: To compare Dierberg's vs. Shop 'n Save in their prices for the 28 items in the shopping basket using the ANOVA t-test.

The first step is to write the null hypothesis and the research hypothesis for these two supermarkets.

For the ANOVA t-test, the null hypothesis is that the population means of the two groups are equal, while the research hypothesis is that the population means of the two groups are not equal (i.e., there is a significant difference between these two means). Since we are comparing Dierberg's (Group 1) vs. Shop 'n Save (Group 3), these hypotheses would be:

$$H_0 : \mu_1 = \mu_3$$

$$H_1 : \mu_1 \neq \mu_3$$

For Group 1 vs. Group 3, the formula for the ANOVA t-test is:

$$ANOVA\ t = \frac{\bar{X}_1 - \bar{X}_2}{s.e.ANOVA} \quad (8.2)$$

where

$$s.e.ANOVA = \sqrt{MS_w \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \quad (8.3)$$

The steps involved in computing this ANOVA t-test are:

1. Find the difference of the sample means for the two groups ($2.44 - 1.69 = 0.75$).
2. Find $1/n_1 + 1/n_3$ (since both groups have 28 supermarket items in them, this becomes: $1/28 + 1/28 = 0.0357 + 0.0357 = 0.0714$)
3. Multiply MS_w times the answer for step 2 ($1.23 \times 0.0714 = 0.0878$)
4. Take the square root of step 3 (SQRT (0.0878) = 0.30)
5. Divide Step 1 by Step 4 to find ANOVA t ($0.75/0.30 = 2.50$)

Note: Since Excel computes all calculations to 16 decimal places, when you use Excel for the above computations, your answer will be 2.54 instead of 2.50 that you will obtain if you use your calculator.

Now, what do we do with this ANOVA t-test result of 2.50? In order to interpret this value of 2.50 correctly, we need to determine the critical value of t for the ANOVA t-test. To do that, we need to find the degrees of freedom for the ANOVA t-test as follows:

8.4.1.1 Finding the Degrees of Freedom for the ANOVA t-Test

Objective: To find the degrees of freedom for the ANOVA t-test.

The degrees of freedom (df) for the ANOVA t-test is found as follows:

df = take the **total sample size of all of the groups** and subtract the number of groups in your study ($n_{\text{TOTAL}} - k$ where k = the number of groups)

In our example, the total sample size of the three groups is 84 since there are 28 prices for each of the three supermarkets, and since there are three groups, $84 - 3$ gives a degrees of freedom for the ANOVA t-test of 81.

If you look up $df = 81$ in the t-table in Appendix E in the **degrees of freedom column (df), which is the second column on the left of this table**, you will find that the critical t-value is 1.96.

Important note: Be sure to use the degrees of freedom column (df) in Appendix E for the ANOVA t-test critical t value

8.4.1.2 Stating the Decision Rule for the ANOVA t-Test

Objective: To learn the decision rule for the ANOVA t-test

Interpreting the result of the ANOVA t-test follows the same decision rule that we used for both the one-group t-test (see Sect. 4.1.6) and the two-group t-test (see Sect. 5.1.8):

If the absolute value of t is less than the critical value of t , we accept the null hypothesis.

or

If the absolute value of t is greater than the critical value of t , we reject the null hypothesis and accept the research hypothesis.

Since we are using a type of t-test, we need to take the absolute value of t . Since the absolute value of 2.50 is greater than the critical t-value of 1.96, we reject the null hypothesis (that the population means of the two groups are equal) and accept the research hypothesis (that the population means of the two groups are significantly different from one another).

This means that our conclusion, in plain English, is as follows:

The average prices of our market basket of items at Dierberg's were significantly higher than the average prices at Shop 'n Save (\$2.44 vs. \$1.69).

Note that this difference in average prices of \$0.75 might not seem like much, but in practical terms, this means that the average prices at Dierberg's are 44% higher than the average prices at Shop 'n Save. This, clearly, is an important difference in prices from these two supermarkets based on our hypothetical data.

8.4.1.3 Performing an ANOVA t-Test Using Excel Commands

Now, let's do these calculations for the ANOVA t-test using Excel with the file you created earlier in this chapter: SUPER6

A52: Dierberg's vs. Shop 'n Save

A54: $1/n$ of Dierberg's + $1/n$ of Shop 'n Save

A56: s.e. of Dierberg's vs. Shop 'n Save

A58: ANOVA t-test

D54: $=(1/28 + 1/28)$ (no spaces between)

D56: $=\text{SQRT}(D48 * D54)$ (no spaces between)

D58: $=(D40 - D42)/D56$ (no spaces between)

You should now have the following results in these cells when you round off all these figures in the ANOVA t-test to two decimal points:

D54: 0.07

D56: 0.30

D58: 2.54

Save this final result under the file name: SUPER7

Print out the resulting spreadsheet so that it fits onto one page like Fig. 8.5 (Hint: Reduce the Page Layout/Scale to Fit to 75%).

For a more detailed explanation of the ANOVA t -test, see Black (2010).

Important note: You are only allowed to perform an ANOVA t -test comparing the population means of two groups when the F -test produces a significant difference between the population means of all of the groups in your study.

It is improper to do any ANOVA t -test when the value of F is less than the critical value of F . Whenever F is less than the critical F , this means that there was no difference between the population means of the groups, and, therefore, that you cannot test to see if there is a difference between the means of any two groups since this would capitalize on chance differences between these two groups.

8.5 End-of-Chapter Practice Problems

- Suppose that you wanted to compare your company's premium brand of tire (Brand A) against two major competitors' brands (B and C). You have set up a laboratory test of the three types of tires, and you have measured the number of simulated miles driven before the tread length reached a pre-determined amount. The hypothetical results are given in Fig. 8.6. Note that the data are in thousands of miles driven (000), so, for example, 63 is really 63,000 miles driven.

TIRE MILEAGE TEST			
(Data are in thousands of miles)			
	Brand A	Brand B	Brand C
	62	61	65
	61	62	67
	62	63	71
	64	60	66
	61	64	65
		59	64
		62	
		63	
		62	
		63	

Fig. 8.6 Worksheet Data for Chap. 8: Practice Problem #1

- (a) Enter these data on an Excel spreadsheet.
- (b) Perform a *one-way ANOVA test* on these data, and show the resulting ANOVA table *underneath* the input data for the three brands of tires.
- (c) If the F-value in the ANOVA table is significant, create an Excel formula to compute the ANOVA t-test comparing the average for Brand A against Brand C and show the results below the ANOVA table on the spreadsheet (put the standard error and the ANOVA t-test value on separate lines of your spreadsheet, and use two decimal places for each value)
- (d) Print out the resulting spreadsheet so that all of the information fits onto one page
- (e) Save the spreadsheet as: TIRE7

Now, write the answers to the following questions using your Excel printout:

1. What are the null hypothesis and the research hypothesis for the ANOVA F-test?
2. What is MS_b on your Excel printout?
3. What is MS_w on your Excel printout?
4. Compute $F = MS_b/MS_w$ using your calculator.
5. What is the critical value of F on your Excel printout?
6. What is the result of the ANOVA F-test?
7. What is the conclusion of the ANOVA F-test in plain English?
8. If the ANOVA F-test produced a significant difference between the three brands in miles driven, what is the null hypothesis and the research hypothesis for the ANOVA t-test comparing Brand A versus Brand C?
9. What is the mean (average) for Brand A on your Excel printout?
10. What is the mean (average) for Brand C on your Excel printout?
11. What are the degrees of freedom (df) for the ANOVA t-test comparing Brand A versus Brand C?
12. What is the critical t value for this ANOVA t-test in Appendix E for these degrees of freedom?
13. Compute the $s.e._{ANOVA}$ using your calculator.
14. Compute the ANOVA t-test value comparing Brand A versus Brand C using your calculator.
15. What is the result of the ANOVA t-test comparing Brand A versus Brand C?
16. What is the conclusion of the ANOVA t-test comparing Brand A versus Brand C in plain English?

Note that since there are three brands of tires, you need to do three ANOVA t-tests to determine what the significant differences are between the tires. *Since you have just completed the ANOVA t-test comparing Brand A versus Brand C, let's do the ANOVA t-test next comparing Brand A versus Brand B.*

17. State the null hypothesis and the research hypothesis comparing Brand A versus Brand B.
18. What is the mean (average) for Brand A on your Excel printout?

19. What is the mean (average) for Brand B on your Excel printout?
20. What are the degrees of freedom (df) for the ANOVA t-test comparing Brand A versus Brand B?
21. What is the critical t value for this ANOVA t-test in Appendix E for these degrees of freedom?
22. Compute the $s.e._{ANOVA}$ for Brand A versus Brand B using your calculator.
23. Compute the ANOVA t-test value comparing Brand A versus Brand B.
24. What is the result of the ANOVA t-test comparing Brand A versus Brand B?
25. What is the conclusion of the ANOVA t-test comparing Brand A versus Brand B in plain English?

The last ANOVA t-test compares Brand B versus Brand C. Let's do that test below:

26. State the null hypothesis and the research hypothesis comparing Brand B versus Brand C.
 27. What is the mean (average) for Brand B on your Excel printout?
 28. What is the mean (average) for Brand C on your Excel printout?
 29. What are the degrees of freedom (df) for the ANOVA t-test comparing Brand B versus Brand C?
 30. What is the critical t value for this ANOVA t-test in Appendix E for these degrees of freedom?
 31. Compute the $s.e._{ANOVA}$ comparing Brand B versus Brand C using your calculator.
 32. Compute the ANOVA t-test value comparing Brand B versus Brand C with your calculator.
 33. What is the result of the ANOVA t-test comparing Brand B versus Brand C?
 34. What is the conclusion of the ANOVA t-test comparing Brand B versus Brand C in plain English?
 35. What is the summary of the three ANOVA t-tests in plain English?
 36. What recommendation would you make to your company about these three brands of tires based on the results of your analysis? Why would you make that recommendation?
2. In an organization with many different departments in different locations, the retention rate of customers in different departments of the organization is an important factor in the success of the organization. If you measure "retention rate" by the percent of customers for a department at the beginning of the year that were still customers at the end of that year, you can compare the departments in your organization in terms of their customer retention rate. Suppose you decide to take a random sample of locations of your organization, and to record the retention rate for each location during the past year for each of three departments. Note that each department can have a different number of locations in order for ANOVA to be used on the data.

Suppose that your random sample produces the hypothetical data given in Fig. 8.7.

RETENTION RATE (last year in percent)		
PRODUCTION	SALES	ENGINEERING
58	79	80
66	92	89
65	84	87
59	86	88
58	88	86
61	89	84
63	90	80
62	92	82
65	89	86
66	91	83
	82	89
	79	

Fig. 8.7 Worksheet Data for Chap. 8: Practice Problem #2

- (a) Enter these data on an Excel spreadsheet.
- (b) Perform a *one-way ANOVA test* on these data, and show the resulting ANOVA table *underneath* the input data for the three departments. Round off all decimal figures to two decimal places, and center all numbers in the ANOVA table.
- (c) If the F-value in the ANOVA table is significant, create an Excel formula to compute the ANOVA t-test comparing retention rate for PRODUCTION against the retention rate for ENGINEERING, and show the results below the ANOVA table on the spreadsheet (put the standard error and the ANOVA t-test value on separate lines of your spreadsheet, and use two decimal places for each value)
- (d) Print out the resulting spreadsheet so that all of the information fits onto one page
- (e) Save the spreadsheet as: Retention6

Now, write the answers to the following questions using your Excel printout:

1. What are the null hypothesis and the research hypothesis for the ANOVA F-test?
2. What is MS_b on your Excel printout?
3. What is MS_w on your Excel printout?
4. Compute $F = MS_b/MS_w$ using your calculator.
5. What is the critical value of F on your Excel printout?

6. What is the result of the ANOVA F-test?
 7. What is the conclusion of the ANOVA F-test in plain English?
 8. If the ANOVA F-test produced a significant difference between the three departments in their retention rate, what is the null hypothesis and the research hypothesis for the ANOVA t-test comparing PRODUCTION versus ENGINEERING?
 9. What is the mean (average) retention rate for PRODUCTION on your Excel printout?
 10. What is the mean (average) retention rate for ENGINEERING on your Excel printout?
 11. What are the degrees of freedom (df) for the ANOVA t-test comparing PRODUCTION versus ENGINEERING?
 12. What is the critical t value for this ANOVA t-test in Appendix E for these degrees of freedom?
 13. Compute the $s.e._{ANOVA}$ using Excel for PRODUCTION versus ENGINEERING.
 14. Compute the ANOVA t-test value comparing PRODUCTION versus ENGINEERING using Excel.
 15. What is the result of the ANOVA t-test comparing PRODUCTION versus ENGINEERING?
 16. What is the conclusion of the ANOVA t-test comparing PRODUCTION versus ENGINEERING in plain English?
3. Suppose that you have been hired as a consultant by Procter and Gamble to analyze the data from a pilot study involving three recent focus groups who were shown four different television commercials for a new type of Crest toothpaste that have not yet been shown on television. The participants were given a ten-item survey to complete after seeing the commercials, and the hypothetical data from question #8 is given in Fig. 8.8 for the four TV commercials.

ITEM #8: "How believable is this commercial to you?"								
1	2	3	4	5	6	7	8	9
not very believable								very believable
Rating for Focus Groups 1, 2, 3 combined								
Television commercial								
A	B	C	D					
2	3	5	6					
3	4	6	7					
5	5	7	4					
4	2	5	5					
5	6	8	3					
3	1	6	8					
6	4	7	2					
4	3	5	6					
3	7	4	7					
7	6	6	5					
2	5	3	8					
1	3	6	9					
3	4	8	5					
5	2	9	6					
6	3	5	7					

Fig. 8.8 Worksheet Data for Chap. 8: Practice Problem #3

- (a) Enter these data on an Excel spreadsheet.
- (b) Perform a *one-way ANOVA test* on these data, and show the resulting ANOVA table *underneath* the input data for the four types of commercials.
- (c) If the F-value in the ANOVA table is significant, create an Excel formula to compute the ANOVA t-test comparing the average for Commercial B against the average for Commercial D, and show the results below the ANOVA table on the spreadsheet (put the standard error and the ANOVA t-test value on separate lines of your spreadsheet, and use two decimal places for each value)
- (d) Print out the resulting spreadsheet so that all of the information fits onto one page
- (e) Save the spreadsheet as: TV6

Now, write the answers to the following questions using your Excel printout:

1. What are the null hypothesis and the research hypothesis for the ANOVA F-test?
2. What is MS_b on your Excel printout?
3. What is MS_w on your Excel printout?
4. Compute $F = MS_b/MS_w$ using your calculator.
5. What is the critical value of F on your Excel printout?
6. What is the result of the ANOVA F-test?
7. What is the conclusion of the ANOVA F-test in plain English?
8. If the ANOVA F-test produced a significant difference between the four types of TV commercials in their believability, what is the null hypothesis and the research hypothesis for the ANOVA t-test comparing Commercial B versus Commercial D?
9. What is the mean (average) for Commercial B on your Excel printout?
10. What is the mean (average) for Commercial D on your Excel printout?
11. What are the degrees of freedom (df) for the ANOVA t-test comparing Commercial B versus Commercial D?
12. What is the critical t value for this ANOVA t-test in Appendix E for these degrees of freedom?
13. Compute the $s.e._{ANOVA}$ using your calculator for Commercial B versus Commercial D.
14. Compute the ANOVA t-test value comparing Commercial B versus Commercial D using your calculator.
15. What is the result of the ANOVA t-test comparing Commercial B versus Commercial D?
16. What is the conclusion of the ANOVA t-test comparing Commercial B versus Commercial D in plain English?

References

- Black, K. Business Statistics: For Contemporary Decision Making (6th ed.). Hoboken, NJ: John Wiley & Sons, Inc., 2010.
- Weiers, R.M. Introduction to Business Statistics (7th ed.). Mason, OH: South-Western Cengage Learning, 2011.

Appendices

Appendix A: Answers to End-of-Chapter Practice Problems

Chapter 1: Practice Problem #1 Answer (see Fig. A.1)

Survey of new-car features						
Panel of female college students (ages 18-24)						
Question #12: If you were to purchase a new car today, how important to you is the feature that "the car parallel parks itself to the curb" by using a computer?						
1	2	3	4	5	6	7
Not Important						Very Important
		RATING				
		5				
		6		n		17
		4				
		3				
		7		Mean		5.059
		6				
		5				
		7		STDEV		1.713
		6				
		7				
		4		s.e.		0.415
		3				
		1				
		7				
		6				
		4				
		5				

Fig. A.1 Answer to Chap. 1: Practice Problem #1

Chapter 1: Practice Problem #2 Answer (see Fig. A.2)

HUMAN RESOURCES MORALE SURVEY						
Item #21: "Management is doing a good job of keeping employee morale at a high level."						
1	2	3	4	5	6	7
Disagree						Agree
		<u>Rating</u>				
		3				
		6				
		5				
		7		n	23	
		2				
		3				
		6		Mean	4.52	
		5				
		4				
		7		STDEV	1.73	
		6				
		1				
		3		s.e.	0.36	
		2				
		4				
		5				
		6				
		4				
		5				
		3				
		6				
		4				
		7				

Fig. A.2 Answer to Chap. 1: Practice Problem #2

Chapter 1: Practice Problem #3 Answer (see Fig. A.3)

Ford Motor Co.				
Number of defects per day for the Ford Focus				
	Day	No. of defects		
	1	6		
	2	8		
	3	14	n	18
	4	12		
	5	6		
	6	8	Mean	11.944
	7	23		
	8	17		
	9	14	STDEV	4.759
	10	16		
	11	18		
	12	12	s.e.	1.122
	13	13		
	14	15		
	15	8		
	16	6		
	17	9		
	18	10		

Fig. A.3 Answer to Chap. 1: Practice Problem #3

Chapter 2: Practice Problem #1 Answer (see Fig. A.4)

FRAME NUMBERS	Duplicate frame numbers	RANDOM NO.
1	44	0.355
2	33	0.311
3	38	0.305
4	43	0.784
5	13	0.569
6	10	0.365
7	50	0.778
8	1	0.851
9	48	0.156
10	61	0.469
11	4	0.708
12	22	0.905
13	40	0.470
14	37	0.093
15	35	0.225
16	60	0.510
17	59	0.173
18	7	0.776
19	17	0.174
20	30	0.999
21	29	0.830
22	47	0.770
23	20	0.220
24	15	0.150
25	12	0.120
26	11	0.110
27	10	0.100
28	9	0.090
29	8	0.080
30	7	0.070
31	6	0.060
32	5	0.050
33	4	0.040
34	3	0.030
35	2	0.020
36	1	0.010
37	0	0.000
38	0	0.000
39	0	0.000
40	0	0.000
41	0	0.000
42	0	0.000
43	0	0.000
44	0	0.000
45	0	0.000
46	0	0.000
47	0	0.000
48	0	0.000
49	0	0.000
50	0	0.000
51	45	0.961
52	28	0.810
53	24	0.241
54	42	0.888
55	11	0.467
56	56	0.977
57	57	0.610
58	54	0.511
59	9	0.697
60	51	0.884
61	39	0.985
62	53	0.760
63	26	0.163

Fig. A.4 Answer to Chap. 2: Practice Problem #1

Chapter 2: Practice Problem #2 Answer (see Fig. A.5)

Fig. A.5 Answer to
Chap. 2: Practice Problem
#2

FRAME NO.	Duplicate frame no.	Random number
1	45	0.955
2	102	0.804
3	16	0.995
4	8	0.976
5	109	0.221
6	64	0.580
7	37	0.509
8	31	0.208
9	27	0.475
10	76	0.471
11	9	0.952
12	70	0.330
13	13	0.481
14	32	0.754
15	56	0.816
16	46	0.986
17	3	0.692
18	98	0.634
19	10	0.526
20	100	0.825
21	29	0.224
90	101	0.964
91	15	0.901
92	61	0.854
93	90	0.059
94	78	0.451
95	69	0.006
96	93	0.621
97	75	0.764
98	59	0.317
99	2	0.805
100	35	0.984
101	20	0.776
102	73	0.398
103	11	0.747
104	24	0.441
105	82	0.637
106	5	0.152
107	17	0.409
108	34	0.963
109	104	0.072
110	51	0.990
111	6	0.455
112	84	0.508
113	96	0.466
114	67	0.650

Chapter 2: Practice Problem #3 Answer (see Fig. A.6)

<i>Chapter 2: Practice Problem #3 Answer</i>		
FRAME NUMBERS	Duplicate frame numbers	Random number
1	47	0.364
2	68	0.637
3	15	0.217
4	69	0.725
5	67	0.192
6	38	0.577
7	43	0.788
8	50	0.527
9	65	0.040
10	40	0.575
11	57	0.189
12	37	0.648
13	22	0.293
14	3	0.832
15	17	0.819
16	60	0.215
17	5	0.670
18	29	0.112
19	74	0.078
20	72	0.766
21	14	0.972
22	41	0.861
23	53	0.495
24	9	0.004
25	19	0.066
26		0.766
	21	0.949
60	26	0.241
61	36	0.626
62	70	0.044
63	39	0.683
64	2	0.378
65	54	0.030
66	44	0.941
67	25	0.599
68	61	0.118
69	23	0.166
70	27	0.722
71	46	0.747
72	35	0.368
73	11	0.429
74	7	0.299
75	12	0.110
76	30	

Fig. A.6 Answer to Chap. 2: Practice Problem #3

Chapter 3: Practice Problem #1 Answer (see Fig. A.7)

St. Louis Post-Dispatch Phone Survey			
Question #4: "How much would you be willing to pay per week for a 6-month weekday/weekend subscription to the Post-Dispatch?"			
Subscription Price (\$)			
4.15	Null hypothesis:	μ	= \$ 3.80
3.75			
3.80	Research hypothesis:	μ	\neq \$ 3.80
4.10			
3.60	n		22
3.60			
3.65	Mean	\$	3.77
4.40			
3.15	STDEV	\$	0.31
4.00			
3.75	s.e.	\$	0.07
4.00			
3.25	95% confidence interval		
3.75			
3.30	lower limit	\$	3.63
3.75			
3.75	upper limit	\$	3.90
3.65			
4.00	Result:	Since the reference value of \$3.80 is inside the confidence interval, we accept the null hypothesis	
4.10			
3.90	Conclusion:	Past subscribers would be willing to pay \$3.80 per week for a 6-month weekday/weekend subscription to the Post-Dispatch	
3.50			
3.75			

\$3.63	-----	---\$3.77---	\$3.80	-----	\$3.90
lower limit		Mean	Ref. value		upper limit

Fig. A.7 Answer to Chap. 3: Practice Problem #1

Chapter 3: Practice Problem #2 Answer (see Fig. A.8)

HUMAN RESOURCES DEPARTMENT						
MORALE SURVEY OF MANAGERS						
Item #24 How would you rate the quality of leadership shown by top management in this company?						
1	2	3	4	5	6	7
very weak						very strong
		Rating				
		5				
		6		Null hypothesis:	$\mu = 4$	
		3				
		4				
		7		Research hypothesis:	$\mu \neq 4$	
		2				
		3				
		4	n		27	
		2				
		5				
		3	Mean		4.00	
		4				
		2				
		2	STDEV		1.52	
		3				
		6				
		5	s.e.		0.29	
		7				
		4				
		6		95% confidence interval		
		4				
		3		lower limit	3.40	
		4				
		2		upper limit	4.60	
		3				
		5				
		4		---- 3.40 ----- 4.00 ----- 4.60 -----		
				lower limit	Mean and Ref. Value	upper limit
			Result:	Since the reference value of 4.00 is inside the confidence interval, we accept the null hypothesis		
			Conclusion:	Managers rated the quality of leadership shown by top management as neither weak nor strong.		

Fig. A.8 Answer to Chap. 3: Practice Problem #2

Chapter 3: Practice Problem #3 Answer (see Fig. A.9)

		Student Advertising Career Conference						
		Survey						
Item #15:	How likely are you to recommend to other advertising students that they attend next year's AAF Student Advertising Career Conference?							
	1	2	3	4	5	6	7	
	Very Unlikely						Very Likely	
		RATING						
		5		Null Hypothesis: $\mu = 4$				
		6		Research hypothesis: $\mu \neq 4$				
		4						
		7						
AAF4		5						
		6	n			19		
		4						
		3						
		1	Mean			4.84		
		2						
		5						
		6	STDEV			1.71		
		7						
		6						
		7	s.e.			0.39		
		6						
		5						
		3	95% confidence interval					
		4						
					lower limit	4.02		
					upper limit	5.67		
		----- 4 ----- 4.02 ----- 4.84 ----- 5.67 -----						
		Ref Value	lower limit		Mean		upper limit	
Result:	Since the reference value is outside of the confidence interval, we reject the null hypothesis and accept the research hypothesis.							
Conclusion:	Students who attended this year's Student Advertising Career Conference were significantly likely to recommend to other advertising students that they attend next year's AAF Student Advertising Career Conference.							

Fig. A.9 Answer to Chap. 3: Practice Problem #3

Chapter 4: Practice Problem #2 Answer (see Fig. A.11)

BOSTON UNIVERSITY M.S. IN ADVERTISING PROGRAM								
Course: Advertising Management								
Item #12: "How would you rate the instructor's ability to explain advertising concepts clearly?"								
	1	2	3	4	5	6	7	
	Poor					6.05	Excellent	
						Mean		
	RATING							
5	Null hypothesis:						$\mu = 4$	
6	Research hypothesis:						$\mu \neq 4$	
4								
7								
6								
5	n				19			
7								
6								
7	Mean				6.05			
5								
6								
7	STDEV				0.91			
6								
7								
5	s.e.				0.21			
6								
7								
6	critical t				2.101			
7								
	t-test				9.82			
	Result: Since the absolute value of 9.82 is greater than the critical t of 2.101, we reject the null hypothesis and accept the research hypothesis.							
	Conclusion: Students who took Advertising Management this past semester rated the ability of the instructor to explain advertising concepts clearly as significantly positive.							

Fig. A.11 Answer to Chap. 4: Practice Problem #2

Chapter 5: Practice Problem #1 Answer (see Fig. A.13)

Boeing Morale Survey			
Note: A high score indicates high job satisfaction, and a low score indicates low job satisfaction			
Group	n	mean	STDEV
1 Males	241	88.20	4.30
2 Females	202	84.80	5.10
Null hypothesis:	$\mu_1 = \mu_2$		
Research hypothesis:	$\mu_1 \neq \mu_2$		
STDEV1 squared / n1			0.077
STDEV2 squared / n2			0.129
E19 + E21			0.205
s.e.			0.453
critical t			1.96
t-test			7.500
Result:	Since the absolute value of 7.500 is greater than the critical t of 1.96, we reject the null hypothesis and accept the research hypothesis		
Conclusion:	Males had significantly higher job satisfaction scores than females at Boeing last month (88.20 vs. 84.80)		

Fig. A.13 Answer to Chap. 5: Practice Problem #1

Chapter 5: Practice Problem #2 Answer (see Fig. A.14)

UNIVERSITY OF ILLINOIS -- URBANA					
GPA OF MS IN ADVERTISING STUDENTS WHO HAVE COMPLETED ALL ADVERTISING REQUIRED COURSES					
MALES	FEMALES				
2.45	2.83	Group	n	Mean	STDEV
2.53	2.74	1 Males	17	3.15	0.42
2.64	2.86	2 Females	15	3.45	0.37
2.72	3.32				
2.85	3.36				
2.96	3.64	Null hypothesis:	μ_1	=	μ_2
3.01	3.56				
3.11	3.56	Research hypothesis:	μ_1	≠	μ_2
3.24	3.64				
3.35	3.37				
3.36	3.67	(n1 - 1) x STDEV1 squared			2.86
3.38	3.91				
3.21	3.92				
3.52	3.64	(n2 - 1) x STDEV2 squared			1.95
3.64	3.71				
3.75		n1 + n2 - 2			30
3.86					
		1/n1 + 1/n2			0.13
		s.e.			0.14
		critical t			2.042
		t-test			-2.09
Result:	Since the absolute value of - 2.09 is greater than the critical t of 2.042, we reject the null hypothesis and accept the research hypothesis.				
Conclusion:	Female MS in Advertising students who have completed all of the required advertising courses had significantly higher GPAs than male advertising students (3.45 vs. 3.15)				

Fig. A.14 Answer to Chap. 5: Practice Problem #2

Chapter 5: Practice Problem #3 Answer (see Fig. A.15)

American Airlines in-flight meal survey							
Question #10: "How likely are you to purchase an in-flight meal on a future flight?"							
	1	2	3	4	5	6	7
Definitely would not purchase		2.36 Vac	3.23 Bus				Definitely would purchase
Group	n	mean	STDEV				
1 Business	64	3.23	1.04				
2 Vacationers	56	2.36	1.35				
Null hypothesis:			$\mu_1 = \mu_2$				
Research hypothesis:			$\mu_1 \neq \mu_2$				
STDEV1 squared / n1				0.02			
STDEV2 squared / n2				0.03			
E25 + E27				0.05			
s.e.				0.22			
critical t				1.96			
t-test				3.91			
Result:	Since the absolute value of 3.91 is greater than the critical t of 1.96, we reject the null hypothesis and accept the research hypothesis						
Conclusion:	Last month on American Airlines, Vacationers were significantly less likely than Business passengers to indicate that they were planning to purchase an in-flight meal on a future flight (2.36 vs. 3.23)						

Fig. A.15 Answer to Chap. 5: Practice Problem #3

Chapter 6: Practice Problem #1 Answer (see Fig. A.16)

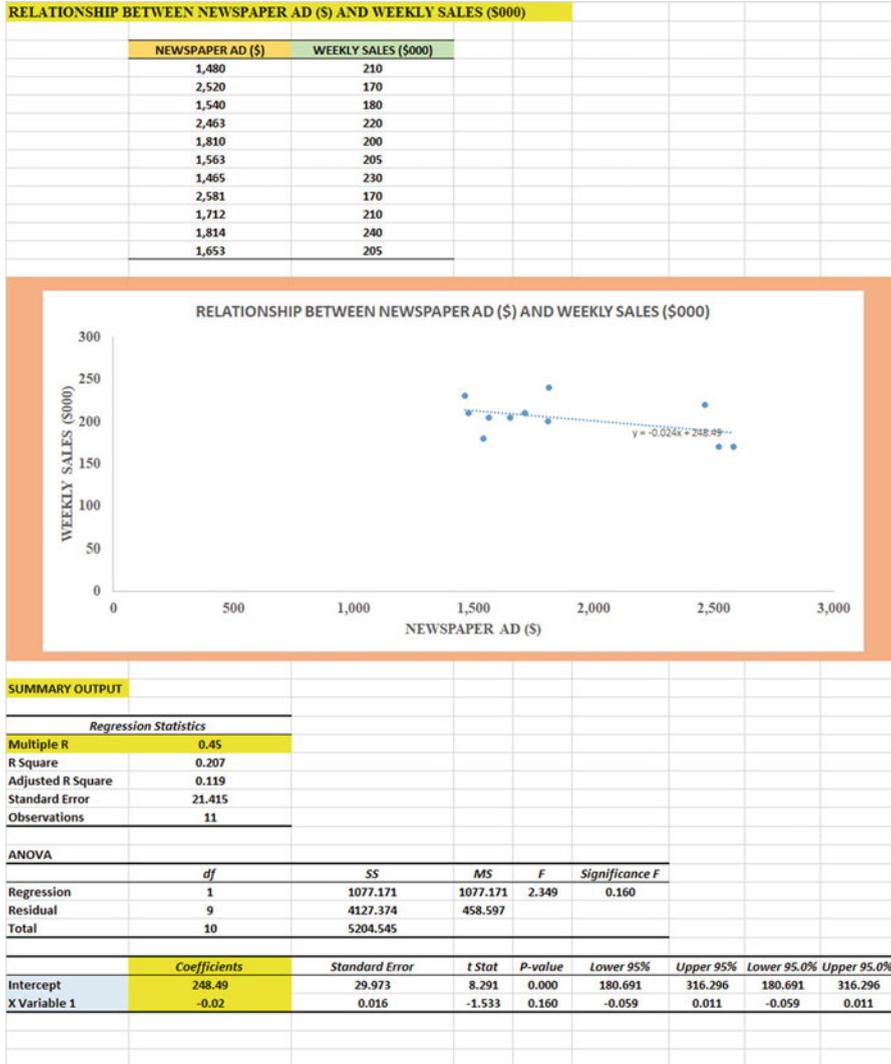


Fig. A.16 Answer to Chap. 6: Practice Problem #1

Chapter 6: Practice Problem #1 (continued)

1. $a = \text{y-intercept} = 248.49$
 $b = \text{slope} = -0.02$ (Note that the slope is negative!)
2. $Y = a + b X$
 $Y = 248.49 - 0.02 X$
3. $r = -.45$ (Note that the correlation is negative!)
4. $Y = 248.49 - 0.02 (2000)$
 $Y = 248.49 - 40$
 $Y = 208.49$
 $Y = \$208,490$
5. About \$220,000

Chapter 6: Practice Problem #2 Answer (see Fig. A.17)

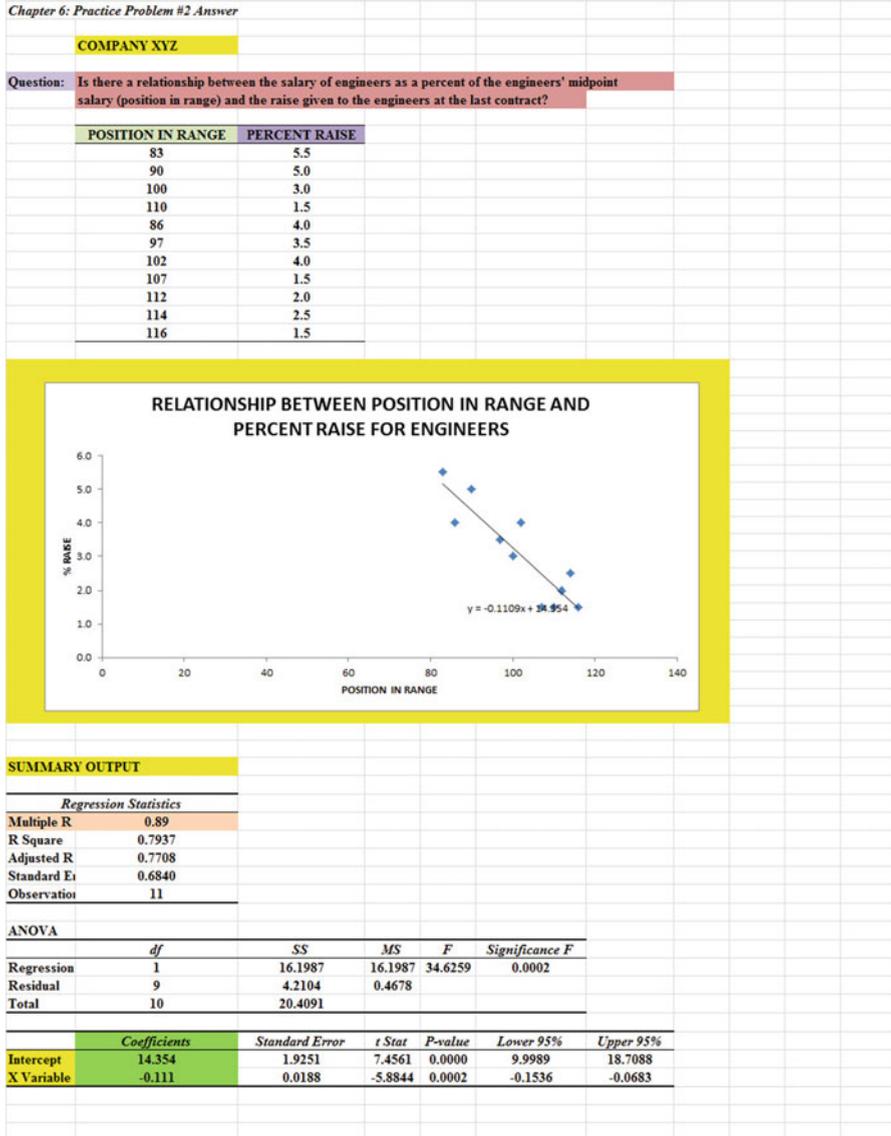


Fig. A.17 Answer to Chap. 6: Practice Problem #2

Chapter 6: Practice Problem #2 (continued)

- (d) $a = y\text{-intercept} = 14.354$
 $b = \text{slope} = -0.111$ (note the minus sign as the slope is negative)
- (e) $Y = a + b X$
 $Y = 14.354 - 0.111 X$
- (f) $r = -.89$ (note the negative correlation!)
- (g) $Y = 14.354 - 0.111 (90)$
 $Y = 14.354 - 9.99$
 $Y = 4.4\%$
- (h) About 2%

Chapter 6: Practice Problem #3 Answer (see Fig. A.18)

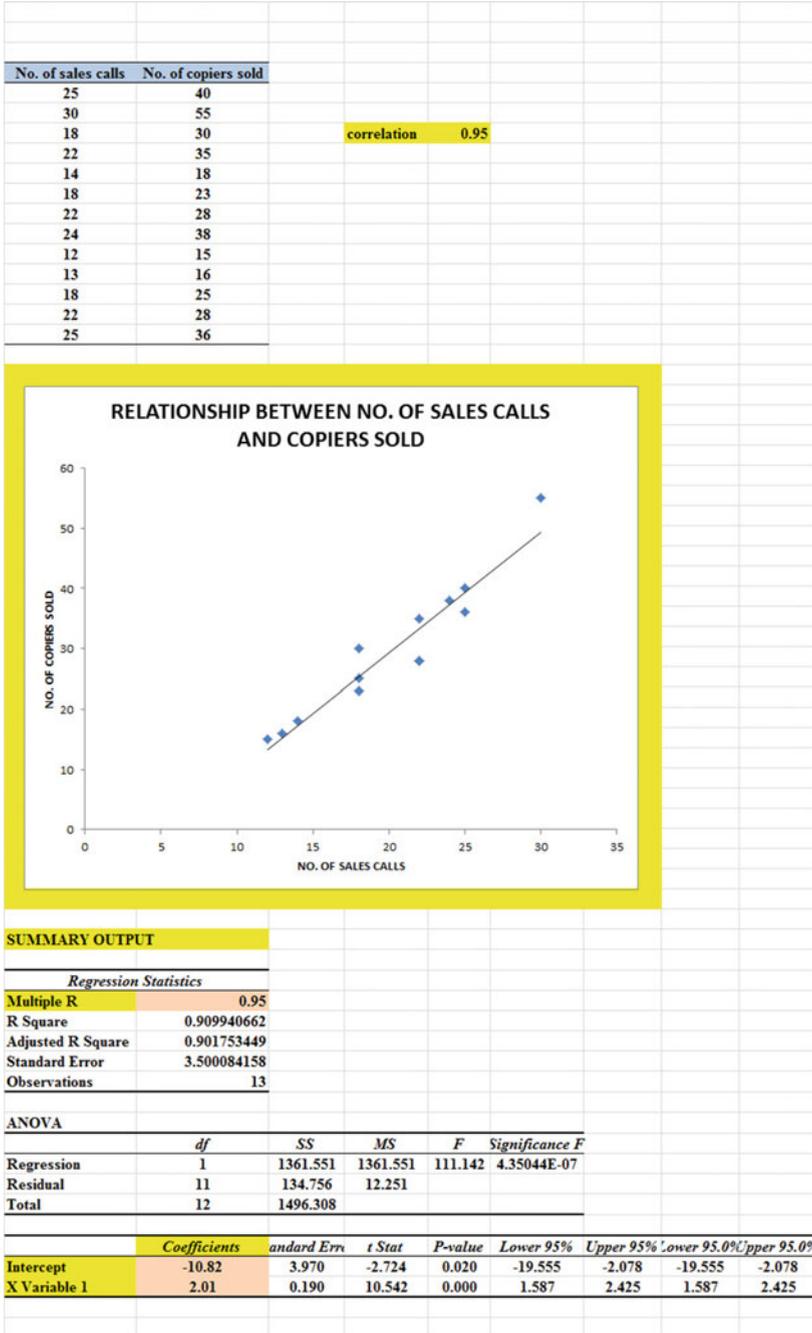


Fig. A.18 Answer to Chap. 6: Practice Problem #3

Chapter 6: Practice Problem #3 (continued)

1. $r = .95$
2. $a = \text{y-intercept} = -10.82$
3. $b = \text{slope} = 2.01$
4. $Y = a + b X$
 $Y = -10.82 + 2.01 X$
5. $Y = -10.82 + 2.01 (25)$
 $Y = -10.82 + 50.25$
 $Y = 39.43$
 $Y = 39 \text{ copiers sold per month}$

Chapter 7: Practice Problem #1 Answer (see Fig. A.19)

Chapter 7: Practice Problem #1 Answer					
GRADUATE RECORD EXAMINATIONS (GRE)					
How well does the GRE predict first-year GPA in an MBA program?					
FIRST-YEAR GPA	GRE VERBAL	GRE QUANTITATIVE	GRE WRITING	UNDERGRAD GPA	
3.25	160	161	5	3.40	
3.42	156	158	4	3.15	
2.85	156	157	2	3.05	
2.65	154	153	1	2.55	
3.65	166	166	6	3.25	
3.16	159	160	3	3.20	
3.56	166	163	4	3.66	
2.35	155	154	2	2.55	
2.86	153	154	3	2.85	
2.95	158	157	4	2.80	
3.15	158	159	4	3.05	
3.45	160	160	5	3.44	
SUMMARY OUTPUT					
Regression Statistics					
Multiple R	0.94				
R Square	0.8825				
Adjusted R Square	0.8154				
Standard Error	0.1676				
Observations	12				
ANOVA					
	df	SS	MS	F	Significance F
Regression	4	1.4777	0.3694	13.1467	0.0023
Residual	7	0.1967	0.0281		
Total	11	1.6744			
	Coefficients	Standard Error	t Stat	P-value	Lower 95%
Intercept	-3.241	4.3231	-0.7496	0.4779	-13.4632
GRE VERBAL	-0.018	0.0388	-0.4590	0.6601	-0.1094
GRE QUANTITATIVE	0.046	0.0561	0.8237	0.4373	-0.0865
GRE WRITING	0.076	0.0654	1.1589	0.2845	-0.0789
UNDERGRAD GPA	0.510	0.2642	1.9303	0.0949	-0.1147
	FIRST-YEAR GPA	GRE VERBAL	GRE QUANTITATIVE	GRE WRITING	UNDERGRAD GPA
FIRST-YEAR GPA	1				
GRE VERBAL	0.79	1			
GRE QUANTITATIVE	0.88	0.94	1		
GRE WRITING	0.83	0.72	0.83	1	
UNDERGRAD GPA	0.88	0.77	0.83	0.70	1

Fig. A.19 Answer to Chap. 7: Practice Problem #1

Chapter 7: Practice Problem #1 (continued)

1. Multiple correlation = $R_{xy} = .94$
2. y-intercept = $a = -3.241$
3. b_1 coefficient = -0.018
4. b_2 coefficient = 0.046
5. b_3 coefficient = 0.076
6. b_4 coefficient = 0.510
7. $Y = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 X_4$
 $Y = -3.241 - 0.018 X_1 + 0.046 X_2 + 0.076 X_3 + 0.510 X_4$
8. $Y = -3.241 - 0.018 (159) + 0.046 (154) + 0.076 (4) + 0.510 (3.05)$
 $Y = -3.241 - 2.862 + 7.084 + 0.304 + 1.556$
 $Y = 8.944 - 6.103$
 $Y = 2.84$
9. 0.88
10. 0.77
11. 0.83
12. 0.70
13. The best predictor of FIRST-YEAR GPA is a tie between GRE QUANTITATIVE and UNDERGRAD GPA ($r = .88$)
14. The four predictors combined predict FIRST-YEAR GPA much better ($R_{xy} = .94$) than the best single predictors by themselves ($r = .88$).

Chapter 7: Practice Problem #2 Answer (see Fig. A.20)

GRADUATE MANAGEMENT ADMISSION TEST (GMAT)					
How well does the GMAT predict first-year GPA in an HRM program?					
FIRST-YEAR GPA	VERBAL	QUANTITATIVE	ANALYTICAL WRITING	INTEGRATED REASONING	
3.25	50	45	4.0	4	
3.67	48	48	4.5	6	
2.80	35	51	5.0	5	
3.05	41	50	5.5	4	
3.45	51	49	4.0	3	
3.33	48	45	3.0	7	
2.75	46	51	4.5	8	
2.95	45	48	5.5	5	
2.60	40	51	6.0	6	
3.67	50	50	4.5	4	
3.75	46	48	3.0	7	
3.42	46	46	4.0	6	
3.15	42	48	5.0	7	
3.26	38	49	4.0	5	
2.96	41	51	5.5	4	
SUMMARY OUTPUT					
<i>Regression Statistics</i>					
Multiple R	0.79				
R Square	0.6228				
Adjusted R Square	0.4720				
Standard Error	0.2566				
Observations	15				
<i>ANOVA</i>					
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	4	1.0872	0.2718	4.1283	0.0314
Residual	10	0.6584	0.0658		
Total	14	1.7456			
	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>
Intercept	3.432	2.5029	1.3714	0.2002	-2.1445
VERBAL	0.022	0.0181	1.2382	0.2439	-0.0179
QUANTITATIVE	0.003	0.0447	0.0758	0.9410	-0.0961
ANALYTICAL WRITING	-0.241	0.1069	-2.2596	0.0474	-0.4796
INTEGRATED REASONING	-0.055	0.0503	-1.0846	0.3036	-0.1668
	<i>FIRST-YEAR GPA</i>	<i>VERBAL</i>	<i>QUANTITATIVE</i>	<i>ANALYTICAL WRITING</i>	<i>INTEGRATED REASONING</i>
FIRST-YEAR GPA	1				
VERBAL	0.62	1			
QUANTITATIVE	-0.50	-0.53	1		
ANALYTICAL WRITING	-0.69	-0.51	0.63	1	
INTEGRATED REASONING	-0.09	-0.07	-0.14	-0.26	1

Fig. A.20 Answer to Chap. 7: Practice Problem #2

1. Multiple correlation = $R_{xy} = .79$
2. y-intercept = $a = 3.432$
3. $b_1 = \text{VERBAL} = 0.022$
4. $b_2 = \text{QUANTITATIVE} = 0.003$
5. $b_3 = \text{WRITING} = -0.241$
6. $b_4 = \text{REASONING} = -0.055$
7. $Y = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 X_4$
 $Y = 3.432 + 0.022 X_1 + 0.003 X_2 - 0.241 X_3 - 0.055 X_4$
8. $Y = 3.432 + 0.022 (48) + 0.003 (46) - 0.241 (4.5) - 0.055 (6)$
 $Y = 3.432 + 1.06 + 0.14 - 1.08 - 0.33$
 $Y = 4.63 - 1.41$
 $Y = 3.22$
9. $+0.62$
10. -0.50
11. -0.69
12. -0.09

Chapter 7: Practice Problem #2 (continued)

- 13. -0.53
- 14. +0.63
- 15. -0.26
- 16. -0.14
- 17. The best predictor of FIRST-YEAR GPA was ANALYTICAL WRITING ($r = -.69$). Note that the best predictor is the “highest number,” whether or not it is positive or negative!
- 18. The four predictors combined predict FIRST-YEAR GPA much better ($R_{xy} = .79$) than the best single predictor by itself ($r = -.69$).

Chapter 7: Practice Problem #3 Answer (see Fig. A.21)

HARDWARE STORES RESULTS				
Sales (\$000)	Direct Mail (\$000)	Billboards (\$000)	TV Ads (\$000)	
15.98	1.05	1.65	1.55	
21.59	1.29	1.55	1.97	
25.47	0.96	1.86	1.64	
16.02	1.51	1.66	1.73	
16.58	1.3	1.45	1.61	
20.1	1.06	1.55	1.82	
17.49	0.86	1.42	1.97	
22.53	0.78	1.55	1.61	
23.98	1.2	1.77	1.13	
18.96	0.83	1.17	1.77	
24.64	1.01	1.76	2.04	
19.52	1.34	1.62	1.63	
19.97	0.55	1.95	1.93	
18.86	1.06	1.55	1.63	
22.79	0.9	1.32	2.27	
20.38	1.88	1.15	1.45	
16.51	0.51	1.19	1.93	
18.32	0.53	1.65	1.44	
16.04	1.21	1.45	1.46	
18.91	0.9	1.45	1.66	
20.09	0.97	1.56	1.59	
19.83	1.01	1.45	1.42	
18.42	1.29	1.47	2.05	

SUMMARY OUTPUT	
Regression Statistics	
Multiple R	0.38
R Square	0.146
Adjusted R Square	0.012
Standard Error	2.756
Observations	23

ANOVA				
	df	SS	MS	F
Regression	3	24.742	8.247	1.086
Residual	19	144.289	7.594	
Total	22	169.031		

	Coefficients	Standard Error	t Stat	P-value
Intercept	8.89	7.424	1.197	0.246
Direct Mail (\$000)	0.45	1.929	0.234	0.817
Billboards (\$000)	5.23	2.944	1.776	0.092
TV Ads (\$000)	1.37	2.359	0.579	0.570

	Sales (\$000)	Direct Mail (\$000)	Billboards (\$000)	TV Ads (\$000)
Sales (\$000)	1			
Direct Mail (\$000)	-0.04	1		
Billboards (\$000)	0.36	-0.16	1	
TV Ads (\$000)	0.06	-0.23	-0.14	1

Fig. A.21 Answer to Chap. 7: Practice Problem #3

Chapter 7: Practice Problem #3 (continued)

Let X_1 = Direct Mail, X_2 = Billboards, and X_3 = TV Ads

(1a) Multiple correlation = $+ .38$

(2b) y-intercept = $a = 8.89$

b_1 = Direct Mail = 0.45

b_2 = Billboards = 5.23

b_3 = TV ads = 1.37

(f) $Y = a + b_1 X_1 + b_2 X_2 + b_3 X_3$

$Y = 8.89 + 0.45 X_1 + 5.23X_2 + 1.37 X_3$

(g) $Y = 8.89 + 0.45(1.88) + 5.23(1.15) + 1.37(1.63)$

$Y = 8.89 + 0.85 + 6.01 + 2.23$

$Y = 17.98$

$Y = \$17,980$

(1a) -0.04

(2b) $+0.36$

(3c) $+0.06$

(4d) -0.23

(5e) -0.14

(6f) The best predictor of Sales was Billboards ($r = .36$).

(7g) The three predictors combined predict Sales only slightly better ($R_{xy} = .38$)

Chapter 8: Practice Problem #1 Answer (see Fig. A.22)

Chapter 8: Practice Problem #1 Answer						
TIRE MILEAGE TEST						
(Data are in thousands of miles)						
	Brand A	Brand B	Brand C			
	62	61	65			
	61	62	67			
	62	63	71			
	64	60	66			
	61	64	65			
		59	64			
		62				
		63				
		62				
		63				
Anova: Single Factor						
SUMMARY						
<i>Groups</i>	<i>Count</i>	<i>Sum</i>	<i>Average</i>	<i>Variance</i>		
Brand A	5	310	62.00	1.50		
Brand B	10	619	61.90	2.32		
Brand C	6	398	66.33	6.27		
ANOVA						
<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F crit</i>
Between Groups	83.00	2	41.50	12.83	0.0003	3.55
Within Groups	58.23	18	3.24			
Total	141.24	20				
Brand A vs. Brand C						
$1/5 + 1/6$		0.37				
s.e. ANOVA		1.09				
ANOVA t-test		-3.98				

Fig. A.22 Answer to Chap. 8: Practice Problem #1

Chapter 8: Practice Problem #1 (continued)

1. Null hypothesis: $\mu_A = \mu_B = \mu_C$
 Research hypothesis: $\mu_A \neq \mu_B \neq \mu_C$
2. $MS_b = 41.50$
3. $MS_w = 3.24$
4. $F = 12.81$
5. critical $F = 3.55$
6. Since the F-value of 12.81 is greater than the critical F value of 3.55, we reject the null hypothesis and accept the research hypothesis.
7. There was a significant difference in the number of miles driven between the three brands of tires.

BRAND A vs. BRAND C

8. Null hypothesis: $\mu_A = \mu_C$
 Research hypothesis: $\mu_A \neq \mu_C$
9. 62
10. 66.33
11. degrees of freedom = $21 - 3 = 18$
12. critical $t = 2.101$
13. $s.e._{ANOVA} = \text{SQRT}(MS_w \times \{1/5 + 1/6\}) = \text{SQRT}(3.24 \times \{0.20 + 0.167\}) = \text{SQRT}(1.19) = 1.09$
14. $ANOVA t = (62 - 66.33)/1.09 = -3.97$
15. Since the absolute value of -3.97 is greater than the critical t of 2.101, we reject the null hypothesis and accept the research hypothesis.
16. Brand C was driven significantly more miles than Brand A (66,000 vs. 62,000).

BRAND A vs. BRAND B

17. Null hypothesis: $\mu_A = \mu_B$
 Research hypothesis: $\mu_A \neq \mu_B$
18. 62
19. 61.9
20. degrees of freedom = $21 - 3 = 18$
21. critical $t = 2.101$
22. $s.e._{ANOVA} = \text{SQRT}(MS_w \times \{1/5 + 1/10\}) = \text{SQRT}(3.24 \times \{0.20 + 0.10\}) = \text{SQRT}(0.972) = 0.99$
23. $ANOVA t = (62 - 61.9)/0.99 = 0.10$
24. Since the absolute value of 0.10 is less than the critical t of 2.101, we accept the null hypothesis.
25. There was no difference in the number of miles driven between Brand A and Brand B.

BRAND B vs. BRAND C

26. Null hypothesis: $\mu_B = \mu_C$
 Research hypothesis: $\mu_B \neq \mu_C$

27. 61.90
28. 66.33
29. degrees of freedom = $21 - 3 = 18$
30. critical $t = 2.101$
31. $s.e._{ANOVA} = \text{SQRT}(MS_w \times \{1/10 + 1/6\}) = \text{SQRT}(3.24 \times \{0.10 + 0.167\})$
 $= \text{SQRT}(0.87) = 0.93$
32. ANOVA $t = (61.90 - 66.33)/0.93 = -4.76$
33. Since the absolute value of -4.76 is greater than the critical t of 2.101, we reject the null hypothesis and accept the research hypothesis.
34. Brand C was driven significantly more miles than Brand B (66,000 vs. 62,000).

SUMMARY

35. Brand C was driven significantly more miles than both Brand A and Brand B. There was no difference in the number of miles driven between Brand A and Brand B.
36. Since our company's Brand A was driven significantly less miles than Brand C, we should never claim in our advertising for Brand A that we last more miles than Brand C. Since our Brand A and Brand B were driven the same number of miles, we should never claim that our tires last longer than Brand B.

Chapter 8: Practice Problem #2 Answer (see Fig. A.23)

RETENTION RATE (last year in percent)						
PRODUCTION	SALES	ENGINEERING				
58	79	80				
66	92	89				
65	84	87				
59	86	88				
58	88	86				
61	89	84				
63	90	80				
62	92	82				
65	89	86				
66	91	83				
	82	89				
	79					
Anova: Single Factor						
SUMMARY						
Groups	Count	Sum	Average	Variance		
PRODUCTION	10	623	62.30	10.23		
SALES	12	1041	86.75	22.39		
ENGINEERING	11	934	84.91	11.09		
ANOVA						
Source of Variation	SS	df	MS	F	P-value	F crit
Between Groups	3891.29	2	1945.64	129.92	1.68E-15	3.32
Within Groups	449.26	30	14.98			
Total	4340.55	32				
PRODUCTION vs. ENGINEERING						
1/n PRODUCTION + 1/n ENGINEERING		0.19				
s.e. PRODUCTION vs. ENGINEERING		1.69				
ANOVA t-test		-13.37				

Fig. A.23 Answer to Chap. 8: Practice Problem #2

Chapter 8: Practice Problem #2 (continued)

Let $X_1 = \text{PRODUCTION}$, $X_2 = \text{SALES}$, and $X_3 = \text{ENGINEERING}$

1. Null hypothesis: $\mu_1 = \mu_2 = \mu_3$
Research hypothesis: $\mu_1 \neq \mu_2 \neq \mu_3$
2. $MS_b = 1945.64$
3. $MS_w = 14.98$
4. $F = 1945.64/14.98 = 129.88$
5. critical $F = 3.32$
6. Since the F-value of 129.88 is greater than the critical F value of 3.32, we reject the null hypothesis and accept the research hypothesis.
7. There was a significant difference between the three departments in retention rate.
8. Null hypothesis: $\mu_1 = \mu_3$
Research hypothesis: $\mu_1 \neq \mu_3$
9. 62.30
10. 84.91
11. degrees of freedom = $33 - 3 = 30$
12. critical $t = 2.042$
13. $s.e._{ANOVA} = 1.69$
14. ANOVA $t = -13.37$
15. Result: Since the absolute value of -13.37 is greater than the critical t of 2.042, we reject the null hypothesis and accept the research hypothesis.
16. Conclusion: ENGINEERING had a significantly higher retention rate than PRODUCTION (85% vs. 62%).

Chapter 8: Practice Problem #3 Answer (see Fig. A.24)

ITEM #8: "How believable is this commercial to you?"									
1	2	3	4	5	6	7	8	9	
not very believable							very believable		
Rating for Focus Groups 1, 2, 3 combined									
Television commercial									
	A	B	C	D					
2	3	5	6						
3	4	6	7						
5	5	7	4						
4	2	5	5						
5	6	8	3						
3	1	6	8						
6	4	7	2						
4	3	5	6						
3	7	4	7						
7	6	6	5						
2	5	3	8						
1	3	6	9						
3	4	8	5						
5	2	9	6						
6	3	5	7						
Anova: Single Factor									
SUMMARY									
Groups	Count	Sum	Average	Variance					
A	15	59	3.93	2.92					
B	15	58	3.87	2.84					
C	15	90	6.00	2.57					
D	15	88	5.87	3.70					
ANOVA									
Source of Variation	SS	df	MS	F	P-value	F crit			
Between Groups	62.18	3	20.73	6.89	0.0005	2.77			
Within Groups	168.40	56	3.01						
Total	230.58	59							
Commercial B vs. Commercial D									
1/15 + 1/15	0.13								
s.e. ANOVA	0.63								
ANOVA t - test	-3.16								

Fig. A.24 Answer to Chap. 8: Practice Problem #3

Chapter 8: Practice Problem #3 (continued)

1. Null hypothesis: $\mu_A = \mu_B = \mu_C = \mu_D$
 Research hypothesis: $\mu_A \neq \mu_B \neq \mu_C \neq \mu_D$
2. $MS_b = 20.73$
3. $MS_w = 3.01$
4. $F = 6.89$
5. critical $F = 2.77$
6. Since the F-value of 6.89 is greater than the critical F value of 2.77, we reject the null hypothesis and accept the research hypothesis.
7. There was a significant difference in the believability of the four television commercials.
8. Null hypothesis: $\mu_B = \mu_D$
 Research hypothesis: $\mu_B \neq \mu_D$
9. 3.87
10. 5.87
11. degrees of freedom = $60 - 4 = 56$
12. critical $t = 1.96$
13. $s.e._{ANOVA} = \text{SQRT}(MS_w \times \{1/15 + 1/15\}) = \text{SQRT}(3.01 \times \{.067 + .067\}) = \text{SQRT}(0.40) = 0.64$
14. ANOVA $t = (3.87 - 5.87)/0.64 = -3.125$
15. Since the absolute value of -3.125 is greater than the critical t of 1.96, we reject the null hypothesis and accept the research hypothesis.
16. Commercial D was significantly more believable than Commercial B (5.87 vs. 3.87).

Appendix B: Practice Test

Chapter 1: Practice Test

Suppose that you have been asked by the manager of the Webster Groves Subaru dealer in St. Louis to analyze the data from a recent survey of its customers. Subaru of America mails a “SERVICE EXPERIENCE SURVEY” to customers who have recently used the Service Department for their car. Let’s try your Excel skills on Item #10e of this survey (see Fig. B.1).

Item #10e:	"Your overall rating of the quality of work performed on your vehicle."									
	1	2	3	4	5	6	7	8	9	10
Unacceptable										Extraordinary
		RATING								
		8								
		5								
		6								
		5								
		4								
		8								
		7								
		7								
		8								
		6								
		7								
		5								
		4								
		8								
		7								
		5								
		7								
		5								
		7								
		6								

Fig. B.1 Worksheet Data for Chap. 1 Practice Test (Practical Example)

- (a) Create an Excel table for these data, and then use Excel to the right of the table to find the sample size, mean, standard deviation, and standard error of the mean for these data. Label your answers, and round off the mean, standard deviation, and standard error of the mean to two decimal places.
- (b) Save the file as: SUBARU8

Chapter 2: Practice Test

Suppose that you wanted to do a personal interview with a random sample of 12 of your company's 42 salespeople as part of a "company morale survey."

- (a) Set up a spreadsheet of frame numbers for these salespeople with the heading: FRAME NUMBERS
- (b) Then, create a separate column to the right of these frame numbers which duplicates these frame numbers with the title: Duplicate frame numbers.
- (c) Then, create a separate column to the right of these duplicate frame numbers called RAND NO. and use the =*RAND()* function to assign random numbers to all of the frame numbers in the duplicate frame numbers column, and change this column format so that three decimal places appear for each random number.
- (d) Sort the *duplicate frame numbers and random numbers* into a random order.
- (e) Print the result so that the spreadsheet fits onto one page.
- (f) Circle on your printout the I.D. number of the first 12 salespeople that you would interview in your company morale survey.
- (g) Save the file as: RAND15

*Important note: Note that everyone who does this problem will generate a different random order of salesperson ID numbers since Excel assigns a different random number each time the *RAND()* command is used. For this reason, the answer to this problem given in this Excel Guide will have a completely different sequence of random numbers from the random sequence that you generate. This is normal and what is to be expected.*

Chapter 3: Practice Test

Suppose that you have been asked to analyze the data from a flight on Southwest Airlines from St. Louis to Boston. Southwest sent an online customer satisfaction survey to a sample of its frequent fliers the day after the flight and asked them to rate their flight on ten-point scales with 1 = extremely dissatisfied, and 10 = extremely satisfied. The data for Item #2c appear in Fig. B.2.

SOUTHWEST AIRLINES ONLINE SURVEY									
Item #2c:	"Please tell us your overall satisfaction with you gate area experience at the airport (gate agent service, facilities, boarding process, and departure time).								
1	2	3	4	5	6	7	8	9	10
extremely dissatisfied									extremely satisfied
			STL-BOS						
			6						
			3						
			8						
			5						
			9						
			10						
			4						
			7						
			6						
			9						
			8						
			7						
			9						
			10						
			7						
			6						
			8						

Fig. B.2 Worksheet Data for Chap. 3 Practice Test (Practical Example)

- (a) Create an Excel table for these data, and use Excel to the right of the table to find the sample size, mean, standard deviation, and standard error of the mean for these data. Label your answers, and round off the mean, standard deviation, and standard error of the mean to two decimal places in number format.
- (b) By hand, write the null hypothesis and the research hypothesis on your printout.
- (c) Use Excel's *TINV function* to find the 95% confidence interval about the mean for these data. Label your answers. Use two decimal places for the confidence interval figures in number format.
- (d) On your printout, draw a diagram of this 95% confidence interval by hand, including the reference value.
- (e) On your spreadsheet, enter the *result*.
- (f) On your spreadsheet, enter the *conclusion in plain English*.
- (g) Print the data and the results so that your spreadsheet fits onto one page.
- (h) Save the file as: south3

Chapter 4: Practice Test

Suppose that you have been asked by the American Marketing Association to analyze the data from the Summer Educators' conference in San Francisco. In order to check your Excel formulas, you have decided to analyze the data for one of these questions before you analyze the data for the entire survey, one item at a time. The conference used five-point scales with 1 = Definitely Would Not, and 5 = Definitely Would. A random sample of the hypothetical data for this one item is given in Fig. B.3.

American Marketing Association				
Summer Educators' Conference in San Francisco, CA				
Item #3: "How likely are you to recommend the Conference to a friend or colleague?"				
1	2	3	4	5
Definitely would not				Definitely Would
Rating				
4				
5				
3				
4				
2				
5				
4				
5				
3				
5				
4				
5				
3				
2				
1				
4				
5				
4				
5				
3				
5				
5				

Fig. B.3 Worksheet Data for Chap. 4 Practice Test (Practical Example)

- (a) Write the null hypothesis and the research hypothesis on your spreadsheet.
- (b) Create a spreadsheet for these data, and then use Excel to find the sample size, mean, standard deviation, and standard error of the mean to the right of the data set. Use number format (three decimal places) for the mean, standard deviation, and standard error of the mean.

- (c) Type the *critical t* from the t-table in Appendix E onto your spreadsheet, and label it.
- (d) Use Excel to compute the t-test value for these data (use three decimal places) and label it on your spreadsheet.
- (e) Type the *result* on your spreadsheet, and then type the *conclusion in plain English* on your spreadsheet.
- (f) Save the file as: BOS2ANSWER

Chapter 5: Practice Test

Suppose that you work for an insurance company and that you have been asked to analyze the data from a marketing research study in which your company was trying to decide whether to use a male model or a female model in an ad in *Business Week* to announce a new type of life insurance policy that would help to provide income toward retirement.

Since the majority of the subscribers to *The Wall Street Journal* are men, an interesting research question would be the following:

Research question: "Does the gender of the model affect adult men's willingness to learn more about how life insurance can provide income for retirement?"

Suppose that you have shown two groups of adult males (ages 25–44) a mockup of an ad such one group of males saw the ad with a male model, while another group of males saw the identical ad except that it had a female model in the ad. (You randomly assigned these males to one of the two experimental groups.) The two groups were kept separate during the experiment and could not interact with one another.

At the end of a 1-h discussion of the mockup ad, the respondents were asked the question given in Fig. B.4.

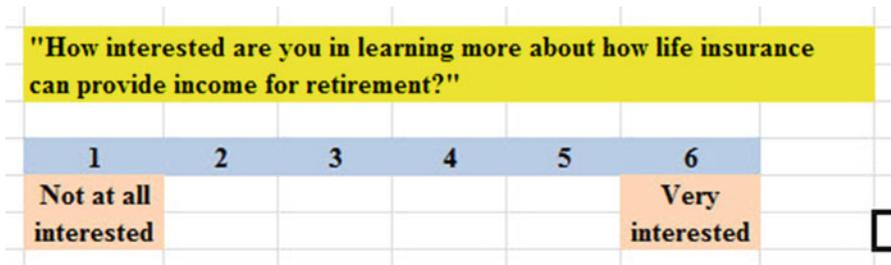


Fig. B.4 Survey Item for a Mockup Ad (Practical Example)

The resulting data for this one item appear in Fig. B.5.

Item: "How interested are you in learning more about how life insurance can provide income for retirement?"					
1	2	3	4	5	6
Not at all interested					Very interested
		Male model	Female model		
		3	4		
		2	6		
		4	5		
		5	3		
		1	4		
		6	6		
		2	6		
		4	5		
		3	3		
		5	5		
		2	4		
		4	3		
		3	5		
		5	4		
		1	6		
		2	5		
		3	5		
		1	6		
		4	4		
		5	6		
		6	3		
		2	4		
		3	6		
		1	5		
		4	6		
		3	4		
		5	4		

Fig. B.5 Worksheet Data for Chap. 5 Practice Test (Practical Example)

- (a) Write the null hypothesis and the research hypothesis.
- (b) Create an Excel table that summarizes these data.
- (c) Use Excel to find the standard error of the difference of the means.
- (d) Use Excel to perform a *two-group t-test*. What is the value of *t* that you obtain (use two decimal places)?
- (e) On your spreadsheet, type the *critical value of t* using the t-table in Appendix E.
- (f) Type the *result* of the test on your spreadsheet.
- (g) Type your *conclusion in plain English* on your spreadsheet.
- (h) Save the file as: Insurance51
- (i) Print the final spreadsheet so that it fits onto one page.

Chapter 6: Practice Test

Suppose that you wanted to study the relationship between DIET (measured in calories allowed per day) and WEIGHT LOSS (measured in kilograms, kg) for adult women between the ages of 30 and 40 who are overweight for their height and body structure, and who all weigh roughly the same number of kilograms before undertaking the weight loss program. You want to test your Excel skills on a random sample of these women based on their weight change over the past 4 months to make sure that you can do this type of research. The hypothetical data appear in Fig. B.6:

Fig. B.6 Worksheet Data for Chap. 6 Practice Test (Practical Example)

RELATIONSHIP BETWEEN DIET AND WEIGHT LOSS	
ADULT WOMEN AGES 30-40	
DIET (calories allowed per day)	WEIGHT LOSS (kg)
900	16.0
1050	12.0
1150	8.0
1275	6.0
1420	3.0
1530	5.5
1610	9.5
1710	2.5
1820	6.0
1875	9.0
1930	6.0
2100	3.0

Create an Excel spreadsheet and enter the data using DIET (calories allowed per day) as the independent variable (predictor) and WEIGHT LOSS (kg) as the dependent variable (criterion). Underneath the table, use Excel's `=correl` function to find the correlation between these two variables. Label the correlation and place it underneath the table; then round off the correlation to two decimal places.

- (a) create an *XY scatterplot* of these two sets of data such that:
- top title: RELATIONSHIP BETWEEN DIET AND WEIGHT LOSS
 - x-axis title: DIET (calories allowed per day)
 - y-axis title: WEIGHT LOSS (kg)
 - move the chart below the table
 - re-size the chart so that it is 8 columns wide and 25 rows long
- (b) Create the *least-squares regression line* for these data on the scatterplot, and add the regression equation to the chart.
- (c) Use Excel to run the regression statistics to find the *equation for the least-squares regression line* for these data and display the results below the chart on your spreadsheet. Use number format (two decimal places) for the correlation and three decimal places for all other decimal figures, including the coefficients.
- (d) Print just the input data and the chart so that this information fits onto one page. Then, print the regression output table on a separate page so that it fits onto that separate page.
- (e) save the file as: DIET3

Answer the following questions using your Excel printout:

1. What is the correlation between DIET and WEIGHT LOSS?
2. What is the y-intercept?
3. What is the slope of the line?
4. What is the regression equation?
5. Use the regression equation to predict the WEIGHT LOSS you would expect for a woman who was practicing a DIET of 1500 calories allowed a day. Show your work on a separate sheet of paper.

Chapter 7: Practice Test

The performance rating given to a manager at an organization is frequently a basis for that manager's promotion opportunities, perceived value to the organization, and, sometimes, even that manager's salary raise. Suppose that you want to study the relationship between the number of years of relevant business experience of a manager, the number of undergraduate or graduate degrees earned by that manager, and that manager's performance rating (rated on a scale where 1 = Poor and 7 = Excellent) at a large, high-tech company. You decide to test your Excel skills on a small sample of mid-level managers at your company to study this relationship.

These hypothetical data appear in Fig. B.7.

Research question:		"Are experience and education good predictors of performance?"	
PERFORMANCE RATING	EXPERIENCE	NO. DEGREES	
7	20	3	
6	15	2	
4	8	2	
1	5	0	
2	6	1	
6	18	3	
5	6	2	
7	10	3	
4	11	2	
5	12	3	
4	8	4	
6	14	3	
5	9	2	

Fig. B.7 Worksheet Data for Chap. 7 Practice Test (Practical Example)

- (a) create an Excel spreadsheet using PERFORMANCE RATING as the criterion, and both the number of years of relevant business experience and the number of undergraduate/graduate degrees earned by the manager as the predictors.
- (b) Save the file as:
Performance2
- (c) Use Excel's *multiple regression* function to find the relationship between these three variables and place the SUMMARY OUTPUT below the table.
- (d) Use number format (two decimal places) for the multiple correlation, and four decimals for the y-intercept, EXPERIENCE, and NO. DEGREES coefficients on the SUMMARY OUTPUT. Use number format (three decimal places) for the other decimal figures in the SUMMARY OUTPUT.
- (e) Print the table and regression results below the table so that they fit onto one page.

Answer the following questions using your Excel printout:

1. What is multiple correlation R_{xy} ?
 2. What is the y-intercept a ?
 3. What is the coefficient for EXPERIENCE b_1 ?
 4. What is the coefficient for NO. DEGREES b_2 ?
 5. What is the multiple regression equation?
 6. Predict the PERFORMANCE RATING you would expect for a manager with 10 years of relevant business experience and three undergraduate/graduate degrees.
- (f) Now, go back to your Excel file and create a correlation matrix for these three variables, and place it underneath the SUMMARY OUTPUT on your spreadsheet.

- (g) Save this file as: Performance3
- (h) Now, print out *just this correlation matrix* on a separate sheet of paper. Answer the following questions using your Excel printout. Be sure to include the plus or minus sign for each correlation:
7. What is the correlation between EXPERIENCE and PERFORMANCE RATING?
 8. What is the correlation between NO. DEGREES and PERFORMANCE RATING?
 9. What is the correlation between EXPERIENCE and NO. DEGREES?
 10. Discuss which of the two predictors is the better predictor of PERFORMANCE RATING.
 11. Explain in words how much better the two predictor variables combined predict PERFORMANCE RATING than the better single predictor by itself.

Chapter 8: Practice Test

Suppose that you worked in R&D for Purina in St. Louis and you were asked to test four flavors of kitten food to see which flavor produces the largest amount of food eaten by kittens. Suppose, further, that the kittens have been matched by age, gender, and species, and randomly assigned to four groups. The resulting amount of food eaten by the kittens appears in the hypothetical data in Fig. B.8. You have been asked to determine if there was a significant difference in the amount of food eaten in these four groups.

FLAVORS OF NEW KITTEN FOOD			
A	B	C	D
12	23	29	38
14	20	27	33
18	17	30	40
11	23	35	34
19	20	33	34
10	28	34	37
17	25	32	43
19	22	35	38
23	28	40	45
16	25	38	39
24			39
15			42

Fig. B.8 Worksheet Data for Chap. 8 Practice Test (Practical Example)

- (a) Enter these data on an Excel spreadsheet.
- (b) On your spreadsheet, write the null hypothesis and the research hypothesis for these data
- (c) Perform a *one-way ANOVA test* on these data, and show the resulting ANOVA table *underneath* the input data for the four types of kitten food.
- (d) If the F-value in the ANOVA table is significant, create an Excel formula to compute the ANOVA t-test comparing the amount of food eaten in Group B against the amount of food eaten in Group D, and show the results below the ANOVA table on the spreadsheet (put the standard error and the ANOVA t-test value on separate lines of your spreadsheet, and use two decimal places for each value)
- (e) Print out the resulting spreadsheet so that all of the information fits onto one page
- (f) On your printout, label by hand the MS (between groups) and the MS (within groups)
- (g) Circle and label the value for F on your printout for the ANOVA of the input data
- (h) Label by hand on the printout the mean for Group B and the mean for Group D that were produced by your ANOVA formulas

Save the spreadsheet as: kitten2

On a separate sheet of paper, now do the following by hand:

- (i) find the critical value of F using the ANOVA Single Factor table that you created
- (j) write a summary of the *results* of the ANOVA test for the input data
- (k) write a summary of the *conclusion* of the ANOVA test in plain English for the input data
- (l) write the null hypothesis and the research hypothesis comparing Group B versus Group D
- (m) compute the degrees of freedom for the *ANOVA t-test* by hand for four flavors.
- (n) write the *critical value of t* for the ANOVA t-test using the table in Appendix E
- (o) write a summary of the *result* of the ANOVA t-test
- (p) write a summary of the *conclusion* of the ANOVA t-test in plain English

Appendix C: Answers to Practice Test

Practice Test Answer: Chapter 1 (see Fig. C.1)

Practice Test Answer: Chapter 1 (see Fig.C.1)										
Item #10e:	"Your overall rating of the quality of work performed on your vehicle."									
	1	2	3	4	5	6	7	8	9	10
Unacceptable										Extraordinary
	RATING									
	8									
	5									
	6					n		20		
	5									
	4									
	8					Mean		6.25		
	7									
	7									
	8					STDEV		1.33		
	6									
	7									
	5					s.e.		0.30		
	4									
	8									
	7									
	5									
	7									
	5									
	7									
	6									

Fig. C.1 Practice Test Answer to Chap. 1 Problem

Practice Test Answer: Chapter 2 (see. Fig. C.2)

FRAME NUMBERS	Duplicate frame numbers	RAND NO.
1	8	0.871
2	22	0.309
3	31	0.658
4	42	0.443
5	4	0.489
6	29	0.370
7	3	0.064
8	21	0.440
9	37	0.026
10	17	0.922
11	34	0.980
12	25	0.930
13	10	0.138
14	41	0.504
15	30	0.884
16	36	0.789
17	13	0.243
18	15	0.250
19	20	0.343
20	14	0.958
21	9	0.779
22	12	0.147
23	38	0.253
24	26	0.476
25	1	0.865
26	5	0.170
27	35	0.410
28	28	0.325
29	24	0.216
30	32	0.439
31	27	0.138
32	19	0.168
33	6	0.326
34	39	0.373
35	2	0.454
36	18	0.777
37	7	0.631
38	11	0.448
39	16	0.412
40	40	0.391
41	33	0.471
42	23	0.865

Fig. C.2 Practice Test Answer to Chap. 2 Problem

Practice Test Answer: Chapter 3 (see. Fig. C.3)

<i>Practice Test Answer: Chapter 3 (see Fig.C.3)</i>			
SOUTHWEST AIRLINES ONLINE SURVEY			
Item #2c:	"Please tell us your overall satisfaction with you gate area experience at the airport (gate agent service, facilities, boarding process, and departure time).		
STL-BOS	Null hypothesis:	$\mu = 5.5$	
6			
3	Research hypothesis:	$\mu \neq 5.5$	
8			
5	n	17	
9			
10			
4	Mean	7.18	
7			
6			
9	STDEV	2.01	
8			
7			
9	s.e.	0.49	
10			
7			
6	95% confidence interval		
8			
	lower limit	6.14	
	upper limit	8.21	
Draw a diagram of the confidence interval			
	----- 5.5 ----- 6.14 ----- 7.18 - ----- 8.21- -----		
	Ref. Value	lower limit	upper limit
Result:	Since the reference value of 5.5 is outside of the confidence interval, we reject the null hypothesis and accept the research hypothesis.		
Conclusion:	Frequent flier passengers on Southwest Airlines flight from St. Louis to Boston were significantly satisfied with their gate experience at the St. Louis airport		

Fig. C.3 Practice Test Answer to Chap. 3 Problem

Practice Test Answer: Chapter 4 (see Fig. C.4)

<i>Practice Test Answer: Chapter 4 (see Fig.C4)</i>			
American Marketing Association			
Summer Educators' Conference in San Francisco, CA			
Item #3: "How likely are you to recommend the Conference to a friend or colleague?"			
Rating	Null hypothesis:	$\mu = 3$	
4			
5	Research hypothesis:	$\mu \neq 3$	
3			
4			
2	n	22	
5			
4			
5	Mean	3.909	
3			
5			
4	STDEV	1.192	
5			
3			
2	s.e.	0.254	
1			
4			
5	critical t	2.080	
4			
5			
3	t-test	3.578	
5			
5			
	Result:	Since the absolute value of 3.578 is greater than the critical t of 2.080, we reject the null hypothesis and accept the research hypothesis.	
	Conclusion:	Attendees at the Summer Educators' Conference of the American Marketing Association in San Francisco were significantly likely to recommend the Conference to a friend or colleague.	

Fig. C.4 Practice Test Answer to Chap. 4 Problem

Practice Test Answer: Chapter 5 (see. Fig. C.5)

Item: "How interested are you in learning more about how life insurance can provide income for retirement?"								
	1	2	3	4	5	6		
	Not at all interested					Very interested		
Male model	Female model		Group	n	Mean	STDEV		
3	4		1 Male model	27	3.30	1.54		
2	6		2 Female model	27	4.70	1.07		
4	5		Null hypothesis:		μ_1	=	μ_2	
5	3		Research hypothesis:		μ_1	≠	μ_2	
1	4							
6	6							
2	6		1/n1 + 1/n2					0.07
4	5							
3	3		(n1 - 1) x STDEV ₁ squared					61.63
5	5							
2	4							
4	3		(n2 - 1) x STDEV ₂ squared					29.63
3	5							
5	4							
1	6		n1 + n2 - 2 (degrees of freedom)					52
2	5							
3	5							
1	6							
4	4		s.e.					0.36
5	6							
6	3		critical t					1.96
2	4							
3	6							
1	5		t-test					-3.90
4	6							
3	4							
5	4							
Result:		Since he absolute value of - 3.90 is greater than the critical t of 1.96, we reject the null hypothesis and accept the research hypothesis.						
Conclusion:		Adult men (ages 25-44) were significantly more interested in learning more about how life insurance can provide income for retirement when a female model was used than when a male model was used in the ad (4.70 vs. 3.30)						

Fig. C.5 Practice Test Answer to Chap. 5 Problem

Practice Test Answer: Chapter 6 (see. Fig. C.6)

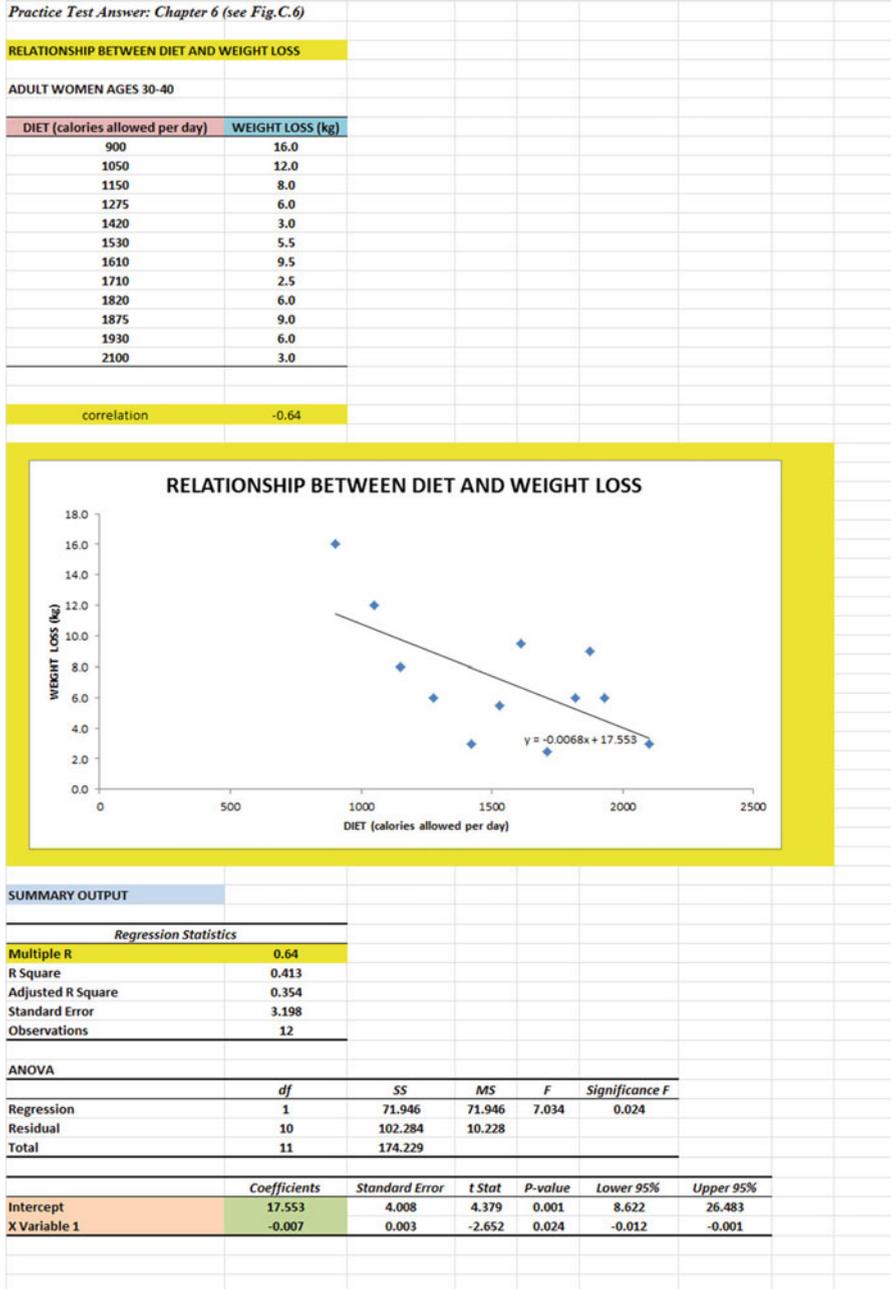


Fig. C.6 Practice Test Answer to Chap. 6 Problem

Practice Test Answer: Chapter 6: (continued)

1. $r = -.64$ (note the negative correlation!)
2. y-intercept = $a = 17.553$
3. slope = $b = -0.007$ (note the negative slope which tells you the correlation is negative!)
4. $Y = a + bX$
 $Y = 17.553 - 0.007X$
5. $Y = 17.553 - 0.007(1500)$
 $Y = 17.553 - 10.5$
 $Y = 7.1$ kg weight loss

Practice Test Answer: Chapter 7 (see Fig. C.7)

Practice Test Answer: Chapter 7 (see Fig.C.7)

PERFORMANCE RATING			EXPERIENCE	NO. DEGREES
7	20	3		
6	15	2		
4	8	2		
1	5	0		
2	6	1		
6	18	3		
5	6	2		
7	10	3		
4	11	2		
5	12	3		
4	8	4		
6	14	3		
5	9	2		

Regression Statistics	
Multiple R	0.84
R Square	0.703
Adjusted R Square	0.644
Standard Error	1.066
Observations	13

ANOVA					
	df	SS	MS	F	Significance F
Regression	2	26.940	13.470	11.850	0.002
Residual	10	11.367	1.137		
Total	12	38.308			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%
Intercept	0.8482	0.858	0.989	0.346	-1.064
EXPERIENCE	0.1916	0.077	2.496	0.032	0.021
NO. DEGREES	0.7922	0.350	2.266	0.047	0.013

	PERFORMANCE RATING	EXPERIENCE	NO. DEGREES
PERFORMANCE RATING	1		
EXPERIENCE	0.74	1	
NO. DEGREES	0.72	0.52	1

Fig. C.7 Practice Test Answer to Chap. 7 Problem

Practice Test Answer: Chapter 7 (continued)

1. Multiple correlation = .84
2. $a = y\text{-intercept} = 0.8482$
3. $b_1 = 0.1916$
4. $b_2 = 0.7922$
5. $Y = a + b_1 X_1 + b_2 X_2$
 $Y = 0.8482 + 0.1916 X_1 + 0.7922 X_2$
6. $Y = 0.8482 + 0.1916 (10) + 0.7922 (3)$
 $Y = 0.8482 + 1.916 + 2.377$
 $Y = 5$
7. +0.74
8. +0.72
9. +0.52
10. The better predictor of PERFORMANCE RATING was EXPERIENCE ($r = .74$).
11. The two predictors combined predicted PERFORMANCE RATING much better at $R_{xy} = .84$.

Practice Test Answer: Chapter 8 (see. Fig. C.8)

				Null hypothesis:	$\mu_A = \mu_B = \mu_C = \mu_D$	
FLAVORS OF NEW KITTEN FOOD						
A	B	C	D	Research hypothesis:	$\mu_A \neq \mu_B \neq \mu_C \neq \mu_D$	
12	23	29	38			
14	20	27	33			
18	17	30	40			
11	23	35	34			
19	20	33	34			
10	28	34	37			
17	25	32	43			
19	22	35	38			
23	28	40	45			
16	25	38	39			
24			39			
15			42			
Anova: Single Factor						
SUMMARY						
Groups	Count	Sum	Average	Variance		
A	12	198	16.50	19.55		
B	10	231	23.10	12.54		
C	10	333	33.30	16.01		
D	12	462	38.50	13.73		
ANOVA						
Source of Variation	SS	df	MS	F	P-value	F crit
Between Groups	3429.55	3	1143.18	73.40	2.58E-16	2.84
Within Groups	623.00	40	15.58			
Total	4052.55	43				
Group B vs. Group D						
1/ n Group B + 1/n Group D		0.18				
s.e. ANOVA		1.69				
ANOVA t-test		-9.11				

Fig. C.8 Practice Test Answer to Chap. 8 Problem

- (f) $MS_b = 1143.18$ and $MS_w = 15.58$
- (g) $F = 73.40$
- (h) Mean Group B = 23.10, and Mean Group D = 38.50
- (i) critical $F = 2.84$
- (j) Results: Since 73.40 is greater than the critical F of 2.84, we reject the null hypothesis and accept the research hypothesis.

Practice Test Answer: Chapter 8 (continued)

- (k) Conclusion: There was a significant difference in the amount of food eaten by the kittens in the four flavors of kitten food
- (l) Null hypothesis: $\mu_B = \mu_D$
Research hypothesis: $\mu_B \neq \mu_D$
- (m) $df = n_{TOTAL} - k = 44 - 4 = 40$
- (n) critical $t = 1.96$
- (o) Result: Since the absolute value of -9.11 is greater than the critical t of 1.96 , we reject the null hypothesis and accept the research hypothesis
- (p) Conclusion: The kittens ate significantly more of Flavor D than Flavor B (38.50 vs. 23.10)

Appendix D: Statistical Formulas

Mean $\bar{X} = \frac{\Sigma X}{n}$

Standard Deviation $STDEV = S = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n - 1}}$

Standard error of the mean $s.e. = S_{\bar{X}} = \frac{S}{\sqrt{n}}$

Confidence interval about the mean $\bar{X} \pm t S_{\bar{X}}$

where $S_{\bar{X}} = \frac{S}{\sqrt{n}}$

One-group t-test $t = \frac{\bar{X} - \mu}{S_{\bar{X}}}$

where $S_{\bar{X}} = \frac{S}{\sqrt{n}}$

Two-group t-test

(a) when both groups have a sample size greater than 30

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}}$$

where $S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$

and where $df = n_1 + n_2 - 2$

(b) when one or both groups have a sample size less than 30

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}}$$

where $S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$

and where $df = n_1 + n_2 - 2$

Correlation

$$r = \frac{\frac{1}{n-1} \Sigma (X - \bar{X})(Y - \bar{Y})}{S_x S_y}$$

where S_x = standard deviation of X

and where S_y = standard deviation of Y

Simple linear regression

$$Y = a + bX$$

where a = y-intercept and b = slope of the line

Multiple regression equation

$$Y = a + b_1X_1 + b_2X_2 + b_3X_3 + \text{etc.}$$

where a = y-intercept

One-way ANOVA F-test

$$F = MS_b / MS_w$$

ANOVA t-test

$$ANOVA \ t = \frac{\bar{X}_1 - \bar{X}_2}{s.e.ANOVA}$$

$$\text{where } s.e.ANOVA = \sqrt{MS_w \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

and where $df = n_{TOTAL} - k$

where $n_{TOTAL} = n_1 + n_2 + n_3 + \text{etc.}$

and where k = the number of groups

Appendix E: t-Table

Critical t-values needed for rejection of the null hypothesis (see Fig. E.1)

Fig. E.1 Critical t-values
Needed for Rejection of the
Null Hypothesis

sample size n	degrees of freedom df	critical t
10	9	2.262
11	10	2.228
12	11	2.201
13	12	2.179
14	13	2.160
15	14	2.145
16	15	2.131
17	16	2.120
18	17	2.110
19	18	2.101
20	19	2.093
21	20	2.086
22	21	2.080
23	22	2.074
24	23	2.069
25	24	2.064
26	25	2.060
27	26	2.056
28	27	2.052
29	28	2.048
30	29	2.045
31	30	2.042
32	31	2.040
33	32	2.037
34	33	2.035
35	34	2.032
36	35	2.030
37	36	2.028
38	37	2.026
39	38	2.024
40	39	2.023
infinity	infinity	1.960

Index

A

- Absolute value of a number, 68, 69
- Analysis of Variance (ANOVA)
 - ANOVA t-test formula (8.2), 175–180
 - degrees of freedom, 176–177, 181, 182, 184, 186, 216, 217, 219, 221, 233
 - Excel commands, 178
 - formula (8.1), 171–173
 - interpreting the Summary Table, 177
 - s.e. formula for ANOVA t-test (8.3), 176
- ANOVA, *see* Analysis of Variance (ANOVA)
- ANOVA t-test, *see* Analysis of Variance (ANOVA)
- Average function, *see* Mean

C

- Centering information within cells, 6–8
- Chart
 - adding the regression equation, 142–144
 - changing the width and height, 6
 - creating a chart, 121–131
 - drawing the regression line onto the chart, 121–131
 - moving the chart, 130
 - printing the spreadsheet, 145–147
 - reducing the scale, 132
 - scatter chart, 123
 - titles, 124, 125, 127
- Column width (changing), 5, 6, 23, 138, 155
- Confidence interval about the mean
 - drawing a picture, 45
 - formula (3.2), 49
 - lower limit, 38–40, 42, 44, 46, 47, 54, 64
 - 95% confident, 38–48

- upper limit, 38–42, 44, 46, 47, 54, 64
- Correlation
 - formula (6.1), 113
 - negative correlation, 109, 111, 112, 140, 145, 208, 240
 - 9 steps for computing, 114–116
 - positive correlation, 109–111, 116, 120, 145
- CORREL function, *see* Correlation
- COUNT function, 9, 54
- Critical t-value, 60, 177, 246

D

- Data Analysis ToolPak, 134–136
- Data/Sort commands, 26
- Degrees of freedom, 87–90, 92, 100, 176, 177, 181, 182, 184, 186, 216, 217, 219, 221

F

- Fill/Series/Columns commands, 4, 5
 - step value/stop value commands, 5, 22
- Formatting numbers
 - currency format, 15–17
 - decimal format, 11–13, 15–17

H

- Home/Fill/Series commands, 5
- Hypothesis testing
 - decision rule, 64
 - null hypothesis, 50–61, 64, 89, 91
 - rating scale hypotheses, 50–53, 57, 69, 91
 - research hypothesis, 50–54, 56–61, 63
 - 7 steps for hypothesis testing, 53–59

Hypothesis testing (*cont.*)

- stating the conclusion, 54, 56, 58, 59, 71, 88
- stating the result, 54, 55, 71, 88

M

- Mean, 1–20, 37–65, 67–82, 84, 85, 94, 113, 118, 173, 181, 182, 184, 186, 222, 225, 226, 242, 244
 - formula (1.1), 1–2
- Multiple correlation
 - correlation matrix, 160–164, 166, 168, 231, 232
 - Excel commands, 156–163, 231, 232
- Multiple regression
 - correlation matrix, 160–164, 166, 168, 231, 232
 - equation (7.1), (7.2), 153–155
 - Excel commands, 156–163, 231, 232
 - predicting Y, 153

N

- Naming a range of cells, 8–9
- Null hypothesis, *see* Hypothesis testing

O

- One-group t-test for the mean
 - absolute value of a number, 68–69
 - formula (4.1), 67–71
 - hypothesis testing, 67–71
 - s.e. formula (4.2), 67–71
 - 7 steps for hypothesis testing, 67–71

P

- Page Layout/Scale to Fit commands, 30
- Population mean, 37–40, 49, 51, 67, 69, 86, 92, 93, 169, 174–177, 180
- Printing a spreadsheet
 - entire worksheet, 145–147
 - part of the worksheet, 145–147
 - printing a worksheet to fit onto one page, 29–33

R

- RAND(), *see* Random number generator
- Random number generator
 - duplicate frame numbers, 23–28, 34, 35, 223

- frame numbers, 21–28, 34, 35, 223
- sorting duplicate frame numbers, 26–29
- Regression, ix, 109–168, 230, 231, 245
- Regression equation
 - adding it to the chart, 142–144
 - formula (6.3), 140
 - negative correlation, 109, 111, 112, 140, 145, 208, 240
 - predicting Y from x, 153
 - slope, b, 140–143, 206, 208, 210, 240, 245
 - writing the regression equation using the Summary Output, 136–140
 - y-intercept, a, 140–143, 159, 206, 208, 210–212, 214, 231, 240, 241, 245
- Regression line, 121–131, 140–145, 148–151, 230
- Research hypothesis, *see* Hypothesis testing

S

- Sample size, 1–20, 38, 41–43, 46, 49, 54, 61, 63–65, 67, 70, 72, 73, 79, 80, 82–87, 89, 92–95, 99–100, 113, 118, 119, 172, 177, 222, 225, 226, 244
 - COUNT function, 9
- Saving a spreadsheet, 13–14
- Scale to Fit commands, 30, 46
- s.e., *see* Standard error of the mean
- Standard deviation, 1–20, 38, 39, 43, 46, 54, 61, 64, 65, 67, 69, 72, 79–85, 89, 90, 93–95, 99, 105, 106, 118, 222, 225, 226, 244, 245
 - formula (1.2), 2
- Standard error of the mean, 1–20, 38–40, 42, 43, 46, 54, 61, 64, 65, 67, 69, 74, 79, 80, 82, 92, 93, 222, 225, 226, 244
 - formula (1.3), 3
- STDEV, *see* Standard deviation (STDEV)

T

- t-table, *see* Appendix E
- Two-group t-test
 - basic table, 85
 - degrees of freedom, 87–90, 92, 100, 177
 - drawing a picture of the means, 91
 - formula (5.2), 92
 - Formula #1 (5.3), 92–99
 - Formula #2 (5.5), 99–104
 - hypothesis testing, 84–92, 100
 - 9 steps in hypothesis testing, 84–92
 - s.e. formula (5.3), (5.5), 99–104