

2-eBook Bundle!

Statistics FOR DUMMIES®

Statistics FOR DUMMIES



Statistics II FOR DUMMIES

Statistics For Dummies,[®] 2-eBook Bundle

Statistics For Dummies,[®] 2nd Edition

Visit www.dummies.com/cheatsheet/statistics to view this book's cheat sheet.

Table of Contents

[Introduction](#)

[About This Book](#)

[Conventions Used in This Book](#)

[What You're Not to Read](#)

[Foolish Assumptions](#)

[How This Book Is Organized](#)

[Part I: Vital Statistics about Statistics](#)

[Part II: Number-Crunching Basics](#)

[Part III: Distributions and the Central Limit Theorem](#)

[Part IV: Guesstimating and Hypothesizing with Confidence](#)

[Part V: Statistical Studies and the Hunt for a Meaningful Relationship](#)

[Part VI: The Part of Tens](#)

[Icons Used in This Book](#)

[Where to Go from Here](#)

[Part I: Vital Statistics about Statistics](#)

[Chapter 1: Statistics in a Nutshell](#)

[Thriving in a Statistical World](#)

[Designing Appropriate Studies](#)

[Surveys](#)

[Experiments](#)

[Collecting Quality Data](#)

Selecting a good sample
Avoiding bias in your data

Creating Effective Summaries

Descriptive statistics
Charts and graphs

Determining Distributions Performing Proper Analyses

Margin of error and confidence intervals
Hypothesis tests
Correlation, regression, and two-way tables

Drawing Credible Conclusions

Reeling in overstated results
Questioning claims of cause and effect

Becoming a Sleuth, Not a Skeptic

Chapter 2: The Statistics of Everyday Life

Statistics and the Media: More Questions than Answers?

Probing popcorn problems
Venturing into viruses
Comprehending crashes
Mulling malpractice
Belaboring the loss of land
Scrutinizing schools
Studying sports
Banking on business news
Touring the travel news
Surveying sexual stats
Breaking down weather reports
Musing about movies
Highlighting horoscopes

Using Statistics at Work

[Delivering babies — and information](#)

[Posing for pictures](#)

[Poking through pizza data](#)

[Statistics in the office](#)

[Chapter 3: Taking Control: So Many Numbers, So Little Time](#)

[Detecting Errors, Exaggerations, and Just Plain Lies](#)

[Checking the math](#)

[Uncovering misleading statistics](#)

[Looking for lies in all the right places](#)

[Feeling the Impact of Misleading Statistics](#)

[Chapter 4: Tools of the Trade](#)

[Statistics: More than Just Numbers](#)

[Grabbing Some Basic Statistical Jargon](#)

[Data](#)

[Data set](#)

[Variable](#)

[Population](#)

[Sample, random, or otherwise](#)

[Statistic](#)

[Parameter](#)

[Bias](#)

[Mean \(Average\)](#)

[Median](#)

[Standard deviation](#)

[Percentile](#)

[Standard score](#)

[Distribution and normal distribution](#)

[Central Limit Theorem](#)

[z-values](#)

[Experiments](#)

[Surveys \(Polls\)](#)

[Margin of error](#)

[Confidence interval](#)

[Hypothesis testing](#)

[p-values](#)

[Statistical significance](#)

[Correlation versus causation](#)

[Part II: Number-Crunching Basics](#)

[Chapter 5: Means, Medians, and More](#)

[Summing Up Data with Descriptive Statistics](#)

[Crunching Categorical Data: Tables and Percents](#)

[Measuring the Center with Mean and Median](#)

[Averaging out to the mean](#)

[Splitting your data down the median](#)

[Comparing means and medians: Histograms](#)

[Accounting for Variation](#)

[Reporting the standard deviation](#)

[Being out of range](#)

[Examining the Empirical Rule \(68-95-99.7\)](#)

[Measuring Relative Standing with Percentiles](#)

[Calculating percentiles](#)

[Interpreting percentiles](#)

[Gathering a five-number summary](#)

[Exploring interquartile range](#)

[Chapter 6: Getting the Picture: Graphing Categorical Data](#)

[Take Another Little Piece of My Pie Chart](#)

[Tallying personal expenses](#)

[Bringing in a lotto revenue](#)

[Ordering takeout](#)

[Projecting age trends](#)

[Raising the Bar on Bar Graphs](#)

[Tracking transportation expenses](#)

[Making a lotto profit](#)

[Tipping the scales on a bar graph](#)
[Pondering pet peeves](#)

[Chapter 7: Going by the Numbers: Graphing Numerical Data](#)

[Handling Histograms](#)

[Making a histogram](#)
[Interpreting a histogram](#)
[Putting numbers with pictures](#)
[Detecting misleading histograms](#)

[Examining Boxplots](#)

[Making a boxplot](#)
[Interpreting a boxplot](#)

[Tackling Time Charts](#)

[Interpreting time charts](#)
[Understanding variability: Time charts versus histograms](#)
[Spotting misleading time charts](#)

[Part III: Distributions and the Central Limit Theorem](#)

[Chapter 8: Random Variables and the Binomial Distribution](#)

[Defining a Random Variable](#)

[Discrete versus continuous](#)
[Probability distributions](#)
[The mean and variance of a discrete random variable](#)

[Identifying a Binomial](#)

[Checking binomial conditions step by step](#)
[No fixed number of trials](#)
[More than success or failure](#)
[Trials are not independent](#)
[Probability of success \(p\) changes](#)

[Finding Binomial Probabilities Using a Formula](#)

Finding Probabilities Using the Binomial Table

Finding probabilities for specific values of X

Finding probabilities for X greater-than, less-than, or between two values

Checking Out the Mean and Standard Deviation of the Binomial

Chapter 9: The Normal Distribution

Exploring the Basics of the Normal Distribution

Meeting the Standard Normal (Z-) Distribution

Checking out Z

Standardizing from X to Z

Finding probabilities for Z with the Z-table

Finding Probabilities for a Normal Distribution

Finding X When You Know the Percent

Figuring out a percentile for a normal distribution

Translating tricky wording in percentile problems

Normal Approximation to the Binomial

Chapter 10: The t-Distribution

Basics of the t-Distribution

Comparing the t- and Z-distributions

Discovering the effect of variability on t-distributions

Using the t-Table

Finding probabilities with the t-table

Figuring percentiles for the t-distribution

Picking out t^* -values for confidence intervals

Studying Behavior Using the t-Table

Chapter 11: Sampling Distributions and the Central Limit Theorem

[Defining a Sampling Distribution](#)
[The Mean of a Sampling Distribution](#)
[Measuring Standard Error](#)

[Sample size and standard error](#)
[Population standard deviation and standard error](#)

[Looking at the Shape of a Sampling Distribution](#)

[Case 1: The distribution of X is normal](#)
[Case 2: The distribution of X is not normal — enter the Central Limit Theorem](#)

[Finding Probabilities for the Sample Mean](#)
[The Sampling Distribution of the Sample Proportion](#)
[Finding Probabilities for the Sample Proportion](#)

[Part IV: Guesstimating and Hypothesizing with Confidence](#)

[Chapter 12: Leaving Room for a Margin of Error](#)

[Seeing the Importance of That Plus or Minus](#)
[Finding the Margin of Error: A General Formula](#)

[Measuring sample variability](#)
[Calculating margin of error for a sample proportion](#)
[Reporting results](#)
[Calculating margin of error for a sample mean](#)
[Being confident you're right](#)

[Determining the Impact of Sample Size](#)

[Sample size and margin of error](#)
[Bigger isn't always \(that much\) better!](#)
[Keeping margin of error in perspective](#)

[Chapter 13: Confidence Intervals: Making Your Best Guesstimate](#)

[Not All Estimates Are Created Equal](#)
[Linking a Statistic to a Parameter](#)
[Getting with the Jargon](#)
[Interpreting Results with Confidence](#)

[Zooming In on Width](#)

[Choosing a Confidence Level](#)

[Factoring In the Sample Size](#)

[Counting On Population Variability](#)

[Calculating a Confidence Interval for a Population Mean](#)

[Case 1: Population standard deviation is known](#)

[Case 2: Population standard deviation is unknown and/or n is small](#)

[Figuring Out What Sample Size You Need](#)

[Determining the Confidence Interval for One Population Proportion](#)

[Creating a Confidence Interval for the Difference of Two Means](#)

[Case 1: Population standard deviations are known](#)

[Case 2: Population standard deviations are unknown and/or sample sizes are small](#)

[Estimating the Difference of Two Proportions](#)

[Spotting Misleading Confidence Intervals](#)

[Chapter 14: Claims, Tests, and Conclusions](#)

[Setting Up the Hypotheses](#)

[Defining the null](#)

[What's the alternative?](#)

[Gathering Good Evidence \(Data\)](#)

[Compiling the Evidence: The Test Statistic](#)

[Gathering sample statistics](#)

[Measuring variability using standard errors](#)

[Understanding standard scores](#)

[Calculating and interpreting the test statistic](#)

[Weighing the Evidence and Making Decisions: p-Values](#)

[Connecting test statistics and p-values](#)

[Defining a p-value](#)

[Calculating a p-value](#)

[Making Conclusions](#)

Setting boundaries for rejecting Ho
Testing varicose veins

Assessing the Chance of a Wrong Decision

Making a false alarm: Type-1 errors
Missing out on a detection: Type-2 errors

Chapter 15: Commonly Used Hypothesis Tests: Formulas and Examples

Testing One Population Mean

Handling Small Samples and Unknown Standard Deviations: The t-Test

Putting the t-test to work

Relating t to Z

Handling negative t-values

Examining the not-equal-to alternative

Testing One Population Proportion

Comparing Two (Independent) Population Averages

Testing for an Average Difference (The Paired t-Test)

Comparing Two Population Proportions

Part V: Statistical Studies and the Hunt for a Meaningful Relationship

Chapter 16: Polls, Polls, and More Polls

Recognizing the Impact of Polls

Getting to the source

Surveying what's hot

Impacting lives

Behind the Scenes: The Ins and Outs of Surveys

Planning and designing a survey

Selecting the sample

Carrying out a survey

Interpreting results and finding problems

Chapter 17: Experiments: Medical Breakthroughs or Misleading Results?

Boiling Down the Basics of Studies

Looking at the lingo of studies

Observing observational studies

Examining experiments

Designing a Good Experiment

Designing the experiment to make comparisons

Selecting the sample size

Choosing the subjects

Making random assignments

Controlling for confounding variables

Respecting ethical issues

Collecting good data

Analyzing the data properly

Making appropriate conclusions

Making Informed Decisions

Chapter 18: Looking for Links: Correlation and Regression

Picturing a Relationship with a Scatterplot

Making a scatterplot

Interpreting a scatterplot

Quantifying Linear Relationships Using the Correlation

Calculating the correlation

Interpreting the correlation

Examining properties of the correlation

Working with Linear Regression

Figuring out which variable is X and which is Y

Checking the conditions

Calculating the regression line

Interpreting the regression line

Putting it all together with an example: The regression line for the crickets

Making Proper Predictions

Explaining the Relationship: Correlation versus Cause and Effect

Chapter 19: Two-Way Tables and Independence

Organizing a Two-Way Table

Setting up the cells

Figuring the totals

Interpreting Two-Way Tables

Singling out variables with marginal distributions

Examining all groups — a joint distribution

Comparing groups with conditional distributions

Checking Independence and Describing Dependence

Checking for independence

Describing a dependent relationship

Cautiously Interpreting Results

Checking for legitimate cause and effect

Projecting from sample to population

Making prudent predictions

Resisting the urge to jump to conclusions

Part VI: The Part of Tens

Chapter 20: Ten Tips for the Statistically Savvy Sleuth

Pinpoint Misleading Graphs

Pie charts

Bar graphs

Time charts

Histograms

Uncover Biased Data

Search for a Margin of Error

Identify Non-Random Samples

[Sniff Out Missing Sample Sizes](#)
[Detect Misinterpreted Correlations](#)
[Reveal Confounding Variables](#)
[Inspect the Numbers](#)
[Report Selective Reporting](#)
[Expose the Anecdote](#)

[Chapter 21: Ten Surefire Exam Score Boosters](#)

[Know What You Don't Know, and then Do Something about It](#)
[Avoid "Yeah-Yeah" Traps](#)

[Yeah-yeah trap #1](#)
[Yeah-yeah trap #2](#)

[Make Friends with Formulas](#)
[Make an "If-Then-How" Chart](#)
[Figure Out What the Question Is Asking](#)
[Label What You're Given](#)
[Draw a Picture](#)
[Make the Connection and Solve the Problem](#)
[Do the Math — Twice](#)
[Analyze Your Answers](#)

[Appendix: Tables for Reference](#)
[Cheat Sheet](#)

Statistics II For Dummies[®]

Visit www.dummies.com/cheatsheet/statistics2 to view this book's cheat sheet.

Table of Contents

[Introduction](#)

[About This Book](#)
[Conventions Used in This Book](#)
[What You're Not to Read](#)
[Foolish Assumptions](#)
[How This Book Is Organized](#)

[Part I: Tackling Data Analysis and Model-Building Basics](#)

[Part II: Using Different Types of Regression to Make Predictions](#)

[Part III: Analyzing Variance with ANOVA](#)

[Part IV: Building Strong Connections with Chi-Square Tests](#)

[Part V: Nonparametric Statistics: Rebels without a Distribution](#)

[Part VI: The Part of Tens](#)

[Icons Used in This Book](#)

[Where to Go from Here](#)

[Part I: Tackling Data Analysis and Model-Building Basics](#)

[Chapter 1: Beyond Number Crunching: The Art and Science of Data Analysis](#)

[Data Analysis: Looking before You Crunch](#)

[Nothing \(not even a straight line\) lasts forever](#)

[Data snooping isn't cool](#)

[No \(data\) fishing allowed](#)

[Getting the Big Picture: An Overview of Stats II](#)

[Population parameter](#)

[Sample statistic](#)

[Confidence interval](#)

[Hypothesis test](#)

[Analysis of variance \(ANOVA\)](#)

[Multiple comparisons](#)

[Interaction effects](#)

[Correlation](#)

[Linear regression](#)

[Chi-square tests](#)

[Nonparametrics](#)

[Chapter 2: Finding the Right Analysis for the Job](#)

[Categorical versus Quantitative Variables](#)

[Statistics for Categorical Variables](#)

[Estimating a proportion](#)

[Comparing proportions](#)

Looking for relationships between categorical variables
Building models to make predictions

Statistics for Quantitative Variables

Making estimates
Making comparisons
Exploring relationships
Predicting y using x

Avoiding Bias
Measuring Precision with Margin of Error
Knowing Your Limitations

Chapter 3: Reviewing Confidence Intervals and Hypothesis Tests

Estimating Parameters by Using Confidence Intervals

Getting the basics: The general form of a confidence interval
Finding the confidence interval for a population mean
What changes the margin of error?
Interpreting a confidence interval

What's the Hype about Hypothesis Tests?

What H_0 and H_a really represent
Gathering your evidence into a test statistic
Determining strength of evidence with a p-value
False alarms and missed opportunities: Type I and II errors
The power of a hypothesis test

Part II: Using Different Types of Regression to Make Predictions

Chapter 4: Getting in Line with Simple Linear Regression

Exploring Relationships with Scatterplots and Correlations

Using scatterplots to explore relationships
Collating the information by using the correlation coefficient

Building a Simple Linear Regression Model

Finding the best-fitting line to model your data

The y-intercept of the regression line

The slope of the regression line

Making point estimates by using the regression line

No Conclusion Left Behind: Tests and Confidence Intervals for Regression

Scrutinizing the slope

Inspecting the y-intercept

Building confidence intervals for the average response

Making the band with prediction intervals

Checking the Model's Fit (The Data, Not the Clothes!)

Defining the conditions

Finding and exploring the residuals

Using r^2 to measure model fit

Scoping for outliers

Knowing the Limitations of Your Regression Analysis

Avoiding slipping into cause-and-effect mode

Extrapolation: The ultimate no-no

Sometimes you need more than one variable

Chapter 5: Multiple Regression with Two X Variables

Getting to Know the Multiple Regression Model

Discovering the uses of multiple regression

Looking at the general form of the multiple regression model

Stepping through the analysis

Looking at x's and y's

Collecting the Data

Pinpointing Possible Relationships

Making scatterplots

Correlations: Examining the bond

Checking for Multicollinearity

Finding the Best-Fitting Model for Two x Variables

Getting the multiple regression coefficients

Interpreting the coefficients

Testing the coefficients

Predicting y by Using the x Variables

Checking the Fit of the Multiple Regression Model

Noting the conditions

Plotting a plan to check the conditions

Checking the three conditions

Chapter 6: How Can I Miss You If You Won't Leave? Regression Model Selection

Getting a Kick out of Estimating Punt Distance

Brainstorming variables and collecting data

Examining scatterplots and correlations

Just Like Buying Shoes: The Model Looks Nice, But Does It Fit?

Assessing the fit of multiple regression models

Model selection procedures

Chapter 7: Getting Ahead of the Learning Curve with Nonlinear Regression

Anticipating Nonlinear Regression

Starting Out with Scatterplots

Handling Curves in the Road with Polynomials

Bringing back polynomials

Searching for the best polynomial model

Using a second-degree polynomial to pass the quiz

Assessing the fit of a polynomial model

Making predictions

Going Up? Going Down? Go Exponential!

Recollecting exponential models

Searching for the best exponential model

Spreading secrets at an exponential rate

Chapter 8: Yes, No, Maybe So: Making Predictions by Using Logistic Regression

Understanding a Logistic Regression Model

How is logistic regression different from other regressions?

Using an S-curve to estimate probabilities

Interpreting the coefficients of the logistic regression model

The logistic regression model in action

Carrying Out a Logistic Regression Analysis

Running the analysis in Minitab

Finding the coefficients and making the model

Estimating p

Checking the fit of the model

Fitting the movie model

Part III: Analyzing Variance with ANOVA

Chapter 9: Testing Lots of Means? Come On Over to ANOVA!

Comparing Two Means with a t-Test

Evaluating More Means with ANOVA

Spitting seeds: A situation just waiting for ANOVA

Walking through the steps of ANOVA

Checking the Conditions

Verifying independence

Looking for what's normal

Taking note of spread

Setting Up the Hypotheses

Doing the F-Test

Running ANOVA in Minitab

Breaking down the variance into sums of squares

Locating those mean sums of squares

Figuring the F-statistic

Making conclusions from ANOVA

What's next?

Checking the Fit of the ANOVA Model

Chapter 10: Sorting Out the Means with Multiple Comparisons

Following Up after ANOVA

Comparing cellphone minutes: An example

Setting the stage for multiple comparison procedures

Pinpointing Differing Means with Fisher and Tukey

Fishing for differences with Fisher's LSD

Using Fisher's new and improved LSD

Separating the turkeys with Tukey's test

Examining the Output to Determine the Analysis

So Many Other Procedures, So Little Time!

Controlling for baloney with the Bonferroni adjustment

Comparing combinations by using Scheffe's method

Finding out whodunit with Dunnett's test

Staying cool with Student Newman-Keuls

Duncan's multiple range test

Going nonparametric with the Kruskal-Wallis test

Chapter 11: Finding Your Way through Two-Way ANOVA

Setting Up the Two-Way ANOVA Model

Determining the treatments

Stepping through the sums of squares

Understanding Interaction Effects

What is interaction, anyway?

Interacting with interaction plots

Testing the Terms in Two-Way ANOVA

Running the Two-Way ANOVA Table

Interpreting the results: Numbers and graphs

Are Whites Whiter in Hot Water? Two-Way ANOVA Investigates

Chapter 12: Regression and ANOVA: Surprise Relatives!

Seeing Regression through the Eyes of Variation

Spotting variability and finding an “x-planation”

Getting results with regression

Assessing the fit of the regression model

Regression and ANOVA: A Meeting of the Models

Comparing sums of squares

Dividing up the degrees of freedom

Bringing regression to the ANOVA table

Relating the F- and t-statistics: The final frontier

Part IV: Building Strong Connections with Chi-Square Tests

Chapter 13: Forming Associations with Two-Way Tables

Breaking Down a Two-Way Table

Organizing data into a two-way table

Filling in the cell counts

Making marginal totals

Breaking Down the Probabilities

Marginal probabilities

Joint probabilities

Conditional probabilities

Trying To Be Independent

Checking for independence between two categories

Checking for independence between two variables

Demystifying Simpson’s Paradox

Experiencing Simpson’s Paradox

Figuring out why Simpson’s Paradox occurs

Keeping one eye open for Simpson’s Paradox

Chapter 14: Being Independent Enough for the Chi-Square Test

The Chi-square Test for Independence

Collecting and organizing the data

Determining the hypotheses

Figuring expected cell counts

Checking the conditions for the test

Calculating the Chi-square test statistic

Finding your results on the Chi-square table

Drawing your conclusions

Putting the Chi-square to the test

Comparing Two Tests for Comparing Two Proportions

Getting reacquainted with the Z-test for two population proportions

Equating Chi-square tests and Z-tests for a two-by-two table

Chapter 15: Using Chi-Square Tests for Goodness-of-Fit (Your Data, Not Your Jeans)

Finding the Goodness-of-Fit Statistic

What's observed versus what's expected

Calculating the goodness-of-fit statistic

Interpreting the Goodness-of-Fit Statistic Using a Chi-Square

Checking the conditions before you start

The steps of the Chi-square goodness-of-fit test

Part V: Nonparametric Statistics: Rebels without a Distribution

Chapter 16: Going Nonparametric

Arguing for Nonparametric Statistics

No need to fret if conditions aren't met

The median's in the spotlight for a change

So, what's the catch?

Mastering the Basics of Nonparametric Statistics

[Sign](#)

[Rank](#)

[Signed rank](#)

[Rank sum](#)

[Chapter 17: All Signs Point to the Sign Test and Signed Rank Test](#)

[Reading the Signs: The Sign Test](#)

[Testing the median](#)

[Estimating the median](#)

[Testing matched pairs](#)

[Going a Step Further with the Signed Rank Test](#)

[A limitation of the sign test](#)

[Stepping through the signed rank test](#)

[Losing weight with signed ranks](#)

[Chapter 18: Pulling Rank with the Rank Sum Test](#)

[Conducting the Rank Sum Test](#)

[Checking the conditions](#)

[Stepping through the test](#)

[Stepping up the sample size](#)

[Performing a Rank Sum Test: Which Real Estate Agent Sells Homes Faster?](#)

[Checking the conditions for this test](#)

[Testing the hypotheses](#)

[Chapter 19: Do the Kruskal-Wallis and Rank the Sums with the Wilcoxon](#)

[Doing the Kruskal-Wallis Test to Compare More than Two Populations](#)

[Checking the conditions](#)

[Setting up the test](#)

[Conducting the test step by step](#)

[Pinpointing the Differences: The Wilcoxon Rank Sum Test](#)

Pairing off with pairwise comparisons

Carrying out comparison tests to see who's different

Examining the medians to see how they're different

Chapter 20: Pointing Out Correlations with Spearman's Rank

Pickin' On Pearson and His Precious Conditions

Scoring with Spearman's Rank Correlation

Figuring Spearman's rank correlation

Watching Spearman at work: Relating aptitude to performance

Part VI: The Part of Tens

Chapter 21: Ten Common Errors in Statistical Conclusions

Claiming These Statistics Prove . . .

It's Not Technically Statistically Significant, But . . .

Concluding That x Causes y

Assuming the Data Was Normal

Only Reporting "Important" Results

Assuming a Bigger Sample Is Always Better

It's Not Technically Random, But . . .

Assuming That 1,000 Responses Is 1,000 Responses

Of Course the Results Apply to the General Population

Deciding Just to Leave It Out

Chapter 22: Ten Ways to Get Ahead by Knowing Statistics

Asking the Right Questions

Being Skeptical

Collecting and Analyzing Data Correctly

Calling for Help

Retracing Someone Else's Steps

Putting the Pieces Together

Checking Your Answers

Explaining the Output

Making Convincing Recommendations

Establishing Yourself as the Statistics Go-To Guy or Gal

Chapter 23: Ten Cool Jobs That Use Statistics

Pollster

Ornithologist (Bird Watcher)

Sportscaster or Sportswriter

Journalist

Crime Fighter

Medical Professional

Marketing Executive

Lawyer

Stock Broker

Appendix: Reference Tables

Cheat Sheet

<https://vk.com/readinglecture>

Statistics FOR DUMMIES®

Learn to:

- Grasp statistical ideas, techniques, formulas, and calculations
- Interpret and critique graphs and charts, determine probability, and work with confidence intervals
- Critique and analyze data from polls and experiments



Deborah J. Rumsey, PhD

Professor of Statistics, The Ohio State University

Statistics For Dummies,[®] 2nd Edition

Visit www.dummies.com/cheatsheet/statistics to view this book's cheat sheet.

Table of Contents

[Introduction](#)

[About This Book](#)

[Conventions Used in This Book](#)

[What You're Not to Read](#)

[Foolish Assumptions](#)

[How This Book Is Organized](#)

[Part I: Vital Statistics about Statistics](#)

[Part II: Number-Crunching Basics](#)

[Part III: Distributions and the Central Limit Theorem](#)

[Part IV: Guesstimating and Hypothesizing with Confidence](#)

[Part V: Statistical Studies and the Hunt for a Meaningful Relationship](#)

[Part VI: The Part of Tens](#)

[Icons Used in This Book](#)

[Where to Go from Here](#)

[Part I: Vital Statistics about Statistics](#)

[Chapter 1: Statistics in a Nutshell](#)

[Thriving in a Statistical World](#)

[Designing Appropriate Studies](#)

[Surveys](#)

[Experiments](#)

[Collecting Quality Data](#)

[Selecting a good sample](#)

[Avoiding bias in your data](#)

[Creating Effective Summaries](#)

[Descriptive statistics](#)

[Charts and graphs](#)

[Determining Distributions](#)

[Performing Proper Analyses](#)

[Margin of error and confidence intervals](#)

[Hypothesis tests](#)

[Correlation, regression, and two-way tables](#)

[Drawing Credible Conclusions](#)

[Reeling in overstated results](#)

[Questioning claims of cause and effect](#)

[Becoming a Sleuth, Not a Skeptic](#)

[Chapter 2: The Statistics of Everyday Life](#)

[Statistics and the Media: More Questions than Answers?](#)

[Probing popcorn problems](#)

[Venturing into viruses](#)

[Comprehending crashes](#)

[Mulling malpractice](#)

[Belaboring the loss of land](#)

[Scrutinizing schools](#)

[Studying sports](#)

[Banking on business news](#)

[Touring the travel news](#)

[Surveying sexual stats](#)

[Breaking down weather reports](#)

[Musing about movies](#)

[Highlighting horoscopes](#)

[Using Statistics at Work](#)

[Delivering babies — and information](#)

[Posing for pictures](#)

[Poking through pizza data](#)

[Statistics in the office](#)

Chapter 3: Taking Control: So Many Numbers, So Little Time

Detecting Errors, Exaggerations, and Just Plain Lies

Checking the math

Uncovering misleading statistics

Looking for lies in all the right places

Feeling the Impact of Misleading Statistics

Chapter 4: Tools of the Trade

Statistics: More than Just Numbers

Grabbing Some Basic Statistical Jargon

Data

Data set

Variable

Population

Sample, random, or otherwise

Statistic

Parameter

Bias

Mean (Average)

Median

Standard deviation

Percentile

Standard score

Distribution and normal distribution

Central Limit Theorem

z-values

Experiments

Surveys (Polls)

Margin of error

Confidence interval

Hypothesis testing

p-values

Statistical significance

Correlation versus causation

Chapter 5: Means, Medians, and More

Summing Up Data with Descriptive Statistics

Crunching Categorical Data: Tables and Percents

Measuring the Center with Mean and Median

Averaging out to the mean

Splitting your data down the median

Comparing means and medians: Histograms

Accounting for Variation

Reporting the standard deviation

Being out of range

Examining the Empirical Rule (68-95-99.7)

Measuring Relative Standing with Percentiles

Calculating percentiles

Interpreting percentiles

Gathering a five-number summary

Exploring interquartile range

Chapter 6: Getting the Picture: Graphing Categorical Data

Take Another Little Piece of My Pie Chart

Tallying personal expenses

Bringing in a lotto revenue

Ordering takeout

Projecting age trends

Raising the Bar on Bar Graphs

Tracking transportation expenses

Making a lotto profit

Tipping the scales on a bar graph

Pondering pet peeves

Chapter 7: Going by the Numbers: Graphing Numerical Data

Handling Histograms

[Making a histogram](#)

[Interpreting a histogram](#)

[Putting numbers with pictures](#)

[Detecting misleading histograms](#)

[Examining Boxplots](#)

[Making a boxplot](#)

[Interpreting a boxplot](#)

[Tackling Time Charts](#)

[Interpreting time charts](#)

[Understanding variability: Time charts versus histograms](#)

[Spotting misleading time charts](#)

[Part III: Distributions and the Central Limit Theorem](#)

[Chapter 8: Random Variables and the Binomial Distribution](#)

[Defining a Random Variable](#)

[Discrete versus continuous](#)

[Probability distributions](#)

[The mean and variance of a discrete random variable](#)

[Identifying a Binomial](#)

[Checking binomial conditions step by step](#)

[No fixed number of trials](#)

[More than success or failure](#)

[Trials are not independent](#)

[Probability of success \(p\) changes](#)

[Finding Binomial Probabilities Using a Formula](#)

[Finding Probabilities Using the Binomial Table](#)

[Finding probabilities for specific values of X](#)

[Finding probabilities for X greater-than, less-than, or between two values](#)

[Checking Out the Mean and Standard Deviation of the Binomial](#)

Chapter 9: The Normal Distribution

Exploring the Basics of the Normal Distribution

Meeting the Standard Normal (Z-) Distribution

Checking out Z

Standardizing from X to Z

Finding probabilities for Z with the Z-table

Finding Probabilities for a Normal Distribution

Finding X When You Know the Percent

Figuring out a percentile for a normal distribution

Translating tricky wording in percentile problems

Normal Approximation to the Binomial

Chapter 10: The t-Distribution

Basics of the t-Distribution

Comparing the t- and Z-distributions

Discovering the effect of variability on t-distributions

Using the t-Table

Finding probabilities with the t-table

Figuring percentiles for the t-distribution

Picking out t*-values for confidence intervals

Studying Behavior Using the t-Table

Chapter 11: Sampling Distributions and the Central Limit Theorem

Defining a Sampling Distribution

The Mean of a Sampling Distribution

Measuring Standard Error

Sample size and standard error

Population standard deviation and standard error

Looking at the Shape of a Sampling Distribution

[Case 1: The distribution of X is normal](#)

[Case 2: The distribution of X is not normal — enter the Central Limit Theorem](#)

[Finding Probabilities for the Sample Mean](#)

[The Sampling Distribution of the Sample Proportion](#)

[Finding Probabilities for the Sample Proportion](#)

[Part IV: Guesstimating and Hypothesizing with Confidence](#)

[Chapter 12: Leaving Room for a Margin of Error](#)

[Seeing the Importance of That Plus or Minus](#)

[Finding the Margin of Error: A General Formula](#)

[Measuring sample variability](#)

[Calculating margin of error for a sample proportion](#)

[Reporting results](#)

[Calculating margin of error for a sample mean](#)

[Being confident you're right](#)

[Determining the Impact of Sample Size](#)

[Sample size and margin of error](#)

[Bigger isn't always \(that much\) better!](#)

[Keeping margin of error in perspective](#)

[Chapter 13: Confidence Intervals: Making Your Best Guesstimate](#)

[Not All Estimates Are Created Equal](#)

[Linking a Statistic to a Parameter](#)

[Getting with the Jargon](#)

[Interpreting Results with Confidence](#)

[Zooming In on Width](#)

[Choosing a Confidence Level](#)

[Factoring In the Sample Size](#)

[Counting On Population Variability](#)

[Calculating a Confidence Interval for a Population Mean](#)

[Case 1: Population standard deviation is known](#)

[Case 2: Population standard deviation is unknown and/or n is small](#)

Figuring Out What Sample Size You Need

Determining the Confidence Interval for One Population Proportion

Creating a Confidence Interval for the Difference of Two Means

Case 1: Population standard deviations are known

Case 2: Population standard deviations are unknown and/or sample sizes are small

Estimating the Difference of Two Proportions

Spotting Misleading Confidence Intervals

Chapter 14: Claims, Tests, and Conclusions

Setting Up the Hypotheses

Defining the null

What's the alternative?

Gathering Good Evidence (Data)

Compiling the Evidence: The Test Statistic

Gathering sample statistics

Measuring variability using standard errors

Understanding standard scores

Calculating and interpreting the test statistic

Weighing the Evidence and Making Decisions: p-Values

Connecting test statistics and p-values

Defining a p-value

Calculating a p-value

Making Conclusions

Setting boundaries for rejecting H_0

Testing varicose veins

Assessing the Chance of a Wrong Decision

Making a false alarm: Type-1 errors

Missing out on a detection: Type-2 errors

Chapter 15: Commonly Used Hypothesis Tests: Formulas and Examples

Testing One Population Mean

Handling Small Samples and Unknown Standard Deviations: The t-Test

Putting the t-test to work

Relating t to Z

Handling negative t-values

Examining the not-equal-to alternative

Testing One Population Proportion

Comparing Two (Independent) Population Averages

Testing for an Average Difference (The Paired t-Test)

Comparing Two Population Proportions

Part V: Statistical Studies and the Hunt for a Meaningful Relationship

Chapter 16: Polls, Polls, and More Polls

Recognizing the Impact of Polls

Getting to the source

Surveying what's hot

Impacting lives

Behind the Scenes: The Ins and Outs of Surveys

Planning and designing a survey

Selecting the sample

Carrying out a survey

Interpreting results and finding problems

Chapter 17: Experiments: Medical Breakthroughs or Misleading Results?

Boiling Down the Basics of Studies

Looking at the lingo of studies

Observing observational studies

Examining experiments

Designing a Good Experiment

Designing the experiment to make comparisons
Selecting the sample size
Choosing the subjects
Making random assignments
Controlling for confounding variables
Respecting ethical issues
Collecting good data
Analyzing the data properly
Making appropriate conclusions

Making Informed Decisions

Chapter 18: Looking for Links: Correlation and Regression

Picturing a Relationship with a Scatterplot

Making a scatterplot
Interpreting a scatterplot

Quantifying Linear Relationships Using the Correlation

Calculating the correlation
Interpreting the correlation
Examining properties of the correlation

Working with Linear Regression

Figuring out which variable is X and which is Y
Checking the conditions
Calculating the regression line
Interpreting the regression line
Putting it all together with an example: The regression line for the crickets

Making Proper Predictions

Explaining the Relationship: Correlation versus Cause and Effect

Chapter 19: Two-Way Tables and Independence

Organizing a Two-Way Table

Setting up the cells

Figuring the totals

Interpreting Two-Way Tables

Singling out variables with marginal distributions

Examining all groups — a joint distribution

Comparing groups with conditional distributions

Checking Independence and Describing Dependence

Checking for independence

Describing a dependent relationship

Cautiously Interpreting Results

Checking for legitimate cause and effect

Projecting from sample to population

Making prudent predictions

Resisting the urge to jump to conclusions

Part VI: The Part of Tens

Chapter 20: Ten Tips for the Statistically Savvy Sleuth

Pinpoint Misleading Graphs

Pie charts

Bar graphs

Time charts

Histograms

Uncover Biased Data

Search for a Margin of Error

Identify Non-Random Samples

Sniff Out Missing Sample Sizes

Detect Misinterpreted Correlations

Reveal Confounding Variables

Inspect the Numbers

Report Selective Reporting

Expose the Anecdote

Chapter 21: Ten Surefire Exam Score Boosters

Know What You Don't Know, and then Do Something about It Avoid “Yeah-Yeah” Traps

Yeah-yeah trap #1

Yeah-yeah trap #2

Make Friends with Formulas

Make an “If-Then-How” Chart

Figure Out What the Question Is Asking

Label What You’re Given

Draw a Picture

Make the Connection and Solve the Problem

Do the Math — Twice

Analyze Your Answers

Appendix: Tables for Reference

Cheat Sheet

<https://vk.com/readinglecture>

Statistics For Dummies,® 2nd Edition

by Deborah J. Rumsey, PhD



Wiley Publishing, Inc.

Statistics For Dummies,® 2nd Edition

Published by
Wiley Publishing, Inc.
111 River St.
Hoboken, NJ 07030-5774
www.wiley.com

Copyright © 2011 by Wiley Publishing, Inc., Indianapolis, Indiana

Published simultaneously in Canada

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400, fax 978-646-8600. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permissions>.

Trademarks: Wiley, the Wiley Publishing logo, For Dummies, the Dummies Man logo, A Reference for the Rest of Us!, The Dummies Way, Dummies Daily, The Fun and Easy Way, Dummies.com, Making Everything Easier, and related trade dress are trademarks or registered trademarks of John Wiley & Sons, Inc. and/or its affiliates in the United States and other countries, and may not be used without written permission. All other trademarks are the property of their respective owners. Wiley Publishing, Inc., is not associated with any product or vendor mentioned in this book.

Limit of Liability/Disclaimer of Warranty: The publisher and the author make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation warranties of fitness for a particular purpose. No warranty may be created or extended by sales or promotional materials. The advice and strategies contained herein may not be suitable for every situation. This work is sold with the understanding that

the publisher is not engaged in rendering legal, accounting, or other professional services. If professional assistance is required, the services of a competent professional person should be sought. Neither the publisher nor the author shall be liable for damages arising herefrom. The fact that an organization or Website is referred to in this work as a citation and/or a potential source of further information does not mean that the author or the publisher endorses the information the organization or Website may provide or recommendations it may make. Further, readers should be aware that Internet Websites listed in this work may have changed or disappeared between when this work was written and when it is read.

For general information on our other products and services, please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993, or fax 317-572-4002.

For technical support, please visit www.wiley.com/techsupport.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Library of Congress Control Number: 2011921775

ISBN: 978-0-470-91108-2

Manufactured in the United States of America

10 9 8 7 6 5 4 3 2 1



About the Author

Deborah J. Rumsey, PhD, is a Statistics Education Specialist and Auxiliary Professor in the Department of Statistics at The Ohio State University. Dr. Rumsey is a Fellow of the American Statistical Association. She has won the Presidential Teaching Award from Kansas State University and has been inducted into the Wall of Inspiration at her high school alma mater, Burlington High School, in Burlington, Wisconsin. She is also the author of *Statistics II For Dummies*, *Statistics Workbook For Dummies*, *Probability For Dummies*, and *Statistics Essentials For Dummies*. She has published numerous papers and given many professional presentations and workshops on the subject of statistics education. She is the original conference designer of the biennial United States Conference on Teaching Statistics (USCOTS). Her passions include being with her family, camping and bird watching, getting seat time on her Kubota tractor, and cheering the Ohio State Buckeyes on to their next national championship.

<https://vk.com/readinglecture>

Dedication

To my husband Eric: My sun rises and sets with you. To my son Clint: I love you up to the moon and back.

Author's Acknowledgments

My heartfelt thanks to Lindsay Lefevere and Kathy Cox for the opportunity to write *For Dummies* books for Wiley; to my project editors Georgette Beatty, Corbin Collins, and Tere Drenth for their unwavering support and vision; to Marjorie Bond, Monmouth College, for agreeing to be my technical editor (again!); to Paul Stephenson, who also provided technical editing; and to Caitie Copple and Janet Dunn for great copy editing.

Special thanks to Elizabeth Stasny, Joan Garfield, Kythrie Silva, Kit Kilen, Peg Steigerwald, Mike O'Leary, Tony Barkauskas, Ken Berk, and Jim Higgins for inspiration and support along the way; and to my entire family for their steadfast love and encouragement.

Publisher's Acknowledgments

We're proud of this book; please send us your comments through our online registration form located at <http://dummies.custhelp.com>. For other comments, please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993, or fax 317-572-4002.

Some of the people who helped bring this book to market include the following:

Acquisitions, Editorial, and Media Development

Project Editor: Corbin Collins

(*Previous Edition: Tere Drenth*)

Senior Project Editor: Georgette Beatty

Executive Editor: Lindsay Sandman

Copy Editor: Caitlin Copple

(*Previous Edition: Janet S. Dunn, PhD*)

Assistant Editor: David Lutton

Technical Editors: Marjorie E. Bond, Paul L. Stephenson III

Editorial Manager: Michelle Hacker

Editorial Supervisor and Reprint Editor: Carmen Krikorian

Editorial Assistant: Jennette ElNaggar

Cover Photos: © iStockphoto.com/Norebbo

Cartoons: Rich Tennant (www.the5thwave.com)

Composition Services

Project Coordinator: Sheree Montgomery

Layout and Graphics: Carrie A. Cesavice, Corrie Socolovitch

Proofreaders: Dwight Ramsey, Shannon Ramsey

Indexer: Christine Karpeles

Publishing and Editorial for Consumer Dummies

Diane Graves Steele, Vice President and Publisher, Consumer Dummies

Kristin Ferguson-Wagstaffe, Product Development Director, Consumer Dummies

Ensley Eikenburg, Associate Publisher, Travel

Kelly Regan, Editorial Director, Travel

Publishing for Technology Dummies

Andy Cummings, Vice President and Publisher, Dummies Technology/General User

Composition Services

Debbie Stailey, Director of Composition Services

<https://vk.com/readinglecture>

Introduction

You get hit with an incredible amount of statistical information on a daily basis. You know what I'm talking about: charts, graphs, tables, and headlines that talk about the results of the latest poll, survey, experiment, or other scientific study. The purpose of this book is to develop and sharpen your skills in sorting through, analyzing, and evaluating all that info, and to do so in a clear, fun, and pain-free way. You also gain the ability to decipher and make important decisions about statistical results (for example, the results of the latest medical studies), while being ever aware of the ways that people can mislead you with statistics. And you see how to do it right when it's your turn to design the study, collect the data, crunch the numbers, and/or draw the conclusions.

This book is also designed to help those of you out there who are taking an introductory statistics class and can use some back-up. You'll gain a working knowledge of the big ideas of statistics and gather a boatload of tools and tricks of the trade that'll help you get ahead of the curve when you take your exams.

This book is chock-full of real examples from real sources that are relevant to your everyday life — from the latest medical breakthroughs, crime studies, and population trends to the latest U.S. government reports. I even address a survey on the worst cars of the millennium! By reading this book, you'll understand how to collect, display, and analyze data correctly and effectively, and you'll be ready to critically examine and make informed decisions about the latest polls, surveys, experiments, and reports that bombard you every day. You even find out how to use crickets to gauge temperature!

You also get to enjoy poking a little fun at statisticians (who take themselves too seriously at times). After all, with the right skills and knowledge, you don't have to be a statistician to understand introductory statistics.

About This Book

This book departs from traditional statistics texts, references, supplemental books, and study guides in the following ways:

- ✓ It includes practical and intuitive explanations of statistical concepts, ideas, techniques, formulas, and calculations found in an introductory statistics course.
- ✓ It shows you clear and concise step-by-step procedures that explain how you can intuitively work through statistics problems.
- ✓ It includes interesting real-world examples relating to your everyday life and workplace.

- ✓ It gives you upfront and honest answers to your questions like, “What does this really mean?” and “When and how will I ever use this?”

Conventions Used in This Book

You should be aware of three conventions as you make your way through this book:

- ✓ **Definition of sample size (n):** When I refer to the size of a sample, I mean the final number of individuals who participated in and provided information for the study. In other words, n stands for the size of the final data set.
- ✓ **Dual-use of the word *statistics*:** In some situations, I refer to statistics as a subject of study or as a field of research, so the word is a singular noun. For example, “Statistics is really quite an interesting subject.” In other situations, I refer to statistics as the plural of *statistic*, in a numerical sense. For example, “The most common statistics are the mean and the standard deviation.”
- ✓ **Use of the word *data*:** You’re probably unaware of the debate raging amongst statisticians about whether the word *data* should be singular (“data is . . .”) or plural (“data are . . .”). It got so bad that recently one group of statisticians had to develop two different versions of a statistics T-shirt: “Messy Data Happens” and “Messy Data Happen.” At the risk of offending some of my colleagues, I go with the plural version of the word *data* in this book.
- ✓ **Use of the term *standard deviation*:** When I use the term *standard deviation*, I mean s , the sample standard deviation. (When I refer to the population standard deviation, I let you know.)

Here are a few other basic conventions to help you navigate this book:

- ✓ I use *italics* to let you know a new statistical term is appearing on the scene.
- ✓ If you see a **boldfaced** term or phrase in a bulleted list, it’s been designated as a keyword or key phrase.
- ✓ Addresses for Web sites appear in monofont.

What You’re Not to Read

I like to think that you won’t skip anything in this book, but I also know you’re a busy person. So to save time, feel free to skip anything marked with the Technical Stuff icon as well as text in sidebars (the shaded gray boxes that appear throughout the book). These items feature information that’s interesting but not crucial to your basic knowledge of

statistics.

Foolish Assumptions

I don't assume that you've had any previous experience with statistics, other than the fact that you're a member of the general public who gets bombarded every day with statistics in the form of numbers, percents, charts, graphs, "statistically significant" results, "scientific" studies, polls, surveys, experiments, and so on.

What I do assume is that you can do some of the basic mathematical operations and understand some of the basic notation used in algebra, such as the variables x and y , summation signs, taking the square root, squaring a number, and so on. If you need to brush up on your algebra skills, check out *Algebra I For Dummies*, 2nd Edition, by Mary Jane Sterling (Wiley).

I don't want to mislead you: You do encounter formulas in this book, because statistics does involve a bit of number crunching. But don't let that worry you. I take you slowly and carefully through each step of any calculations you need to do. I also provide examples for you to work along with this book, so that you can become familiar and comfortable with the calculations and make them your own.

How This Book Is Organized

This book is organized into five parts that explore the major areas of introductory statistics, along with a final part that offers some quick top-ten nuggets for your information and enjoyment. Each part contains chapters that break down each major area of statistics into understandable pieces.

Part I: Vital Statistics about Statistics

This part helps you become aware of the quantity and quality of statistics you encounter in your workplace and your everyday life. You find out that a great deal of that statistical information is incorrect, either by accident or by design. You take a first step toward becoming statistically savvy by recognizing some of the tools of the trade, developing an overview of statistics as a process for getting and interpreting information, and getting up to speed on some statistical jargon.

Part II: Number-Crunching Basics

This part helps you become more familiar and comfortable with making, interpreting, and

evaluating data displays (otherwise known as charts, graphs, and so on) for different types of data. You also find out how to summarize and explore data by calculating and combining some commonly used statistics as well as some statistics you may not know about yet.

Part III: Distributions and the Central Limit Theorem

In this part, you get into all the details of the three most common statistical distributions: the binomial distribution, the normal (and standard normal, also known as Z-distribution), and the *t*-distribution. You discover the characteristics of each distribution and how to find and interpret probabilities, percentiles, means, and standard deviations. You also find measures of relative standing (like percentiles).

Finally, you discover how statisticians measure variability from sample to sample and why a measure of precision in your sample results is so important. And you get the lowdown on what some statisticians describe as the “Crowning Jewel of all Statistics”: the Central Limit Theorem (CLT). I don’t use quite this level of flourishing language to describe the CLT; I just tell my students it’s an MDR (“Mighty Deep Result”; coined by my PhD adviser). As for how my students describe their feelings about the CLT, I’ll leave that to your imagination.

Part IV: Guesstimating and Hypothesizing with Confidence

This part focuses on the two methods for taking the results from a sample and generalizing them to make conclusions about an entire population. (Statisticians call this process *statistical inference*.) These two methods are confidence intervals and hypothesis tests.

In this part, you use confidence intervals to come up with good estimates for one or two population means or proportions, or for the difference between them (for example, the average number of hours teenagers spend watching TV per week or the percentage of men versus women in the United States who take arthritis medicine every day). You get the nitty-gritty on how confidence intervals are formed, interpreted, and evaluated for correctness and credibility. You explore the factors that influence the width of a confidence interval (such as sample size) and work through formulas, step-by-step calculations, and examples for the most commonly used confidence intervals.

The hypothesis tests in this part show you how to use your data to test someone’s claim about one or two population means or proportions, or the difference between them. (For example, a company claims their packages are delivered in two days on average — is this true?) You discover how researchers (should) go about forming and testing hypotheses and how you can evaluate their results for accuracy and credibility. You also get detailed

step-by-step directions and examples for carrying out and interpreting the results of the most commonly used hypothesis tests.

Part V: Statistical Studies and the Hunt for a Meaningful Relationship

This part gives an overview of surveys, experiments, and observational studies. You find out what these studies do, how they are conducted, what their limitations are, and how to evaluate them to determine whether you should believe the results.

You also get all the details on how to examine pairs of numerical variables and categorical variables to look for relationships; this is the object of a great number of studies. For pairs of categorical variables, you create two-way tables and find joint, conditional, and marginal probabilities and distributions. You check for independence, and if a dependent relationship is found, you describe the nature of the relationship using probabilities. For numerical variables you create scatterplots, find and interpret correlation, perform regression analyses, study the fit of the regression line and the impact of outliers, describe the relationship using the slope, and use the line to make predictions. All in a day's work!

Part VI: The Part of Tens

This quick and easy part shares ten ways to be a statistically savvy sleuth and root out suspicious studies and results, as well as ten surefire ways to boost your statistics exam score.

Some statistical calculations involve the use of statistical tables, and I provide quick and easy access to all the tables you need for this book in the appendix. These tables are the Z-table (for the standard normal, also called the Z-distribution), the t -table (for the t -distribution), and the binomial table (for — you guessed it — the binomial distribution). Instructions and examples for using these three tables are provided in their corresponding sections of this book.

Icons Used in This Book

Icons are used in this book to draw your attention to certain features that occur on a regular basis. Here's what they mean:



This icon refers to helpful hints, ideas, or shortcuts that you can use to save time. It also highlights alternative ways to think about a particular concept.



This icon is reserved for particular ideas that I hope you'll remember long after you read this book.



This icon refers to specific ways that researchers or the media can mislead you with statistics and tells you what you can do about it. It also points out potential problems and cautions to keep an eye out for on exams.



This icon is a sure bet if you have a special interest in understanding the more technical aspects of statistical issues. You can skip this icon if you don't want to get into the gory details.

Where to Go from Here

This book is written in such a way that you can start anywhere and still be able to understand what's going on. So you can take a peek at the table of contents or the index, look up the information that interests you, and flip to the page listed. However if you have a specific topic in mind and are eager to dive into it, here are some directions:

- ✓ To work on finding and interpreting graphs, charts, means or medians, and the like, head to Part II.
- ✓ To find info on the normal, Z-, t-, or binomial distributions or the Central Limit Theorem, see Part III.
- ✓ To focus on confidence intervals and hypothesis tests of all shapes and sizes, flip to Part IV.
- ✓ To delve into surveys, experiments, regression, and two-way tables, see Part V.

Or if you aren't sure where you want to start, you may just go with Chapter 1 for the big picture and then plow your way through the rest of the book. Happy reading!

<https://vk.com/readinglecture>

Part I

Vital Statistics about Statistics

The 5th Wave

By Rich Tennant



"Is it just me or did the whole '50% satisfaction' statistic seem a little unimpressive?"

In this part . . .

When you turn on the TV or open a newspaper, you're bombarded with numbers, charts, graphs, and statistical results. From today's poll to the latest major medical breakthroughs, the numbers just keep coming. Yet much of the statistical information you're asked to consume is actually wrong — by accident or even by design. How is a person to know what to believe? By doing a lot of good detective work.

This part helps awaken the statistical sleuth that lies within you by exploring how statistics affect your everyday life and your job, how bad much of the information out there really is, and what you can do about it. This part also helps you get up to speed with some useful statistical jargon.

Chapter 1

Statistics in a Nutshell

In This Chapter

- ▶ Finding out what the process of statistics is all about
 - ▶ Gaining success with statistics in your everyday life, your career, and in the classroom
-

The world today is overflowing with data to the point where anyone (even me!) can be overwhelmed. I wouldn't blame you if you were cynical right now about statistics you read about in the media — I am too at times. The good news is that while a great deal of misleading and incorrect information is lying out there waiting for you, a lot of great stuff is also being produced; for example, many studies and techniques involving data are helping improve the quality of our lives. Your job is to be able to sort out the good from the bad and be confident in your ability to do that. Through a strong understanding of statistics and statistical procedures, you gain power and confidence with numbers in your everyday life, in your job, and in the classroom. That's what this book is all about.

In this chapter, I give you an overview of the role statistics plays in today's data-packed society and what you can do to not only survive but thrive. You get a much broader view of statistics as a partner in the scientific method — designing effective studies, collecting good data, organizing and analyzing the information, interpreting the results, and making appropriate conclusions. (And you thought statistics was just number-crunching!)

Thriving in a Statistical World

It's hard to get a handle on the flood of statistics that affect your daily life in large and small ways. It begins the moment you wake up in the morning and check the news and listen to the meteorologist give you her predictions for the weather based on her statistical analyses of past data and present weather conditions. You pore over nutritional information on the side of your cereal box while you eat breakfast. At work you pull numbers from charts and tables, enter data into spreadsheets, run diagnostics, take measurements, perform calculations, estimate expenses, make decisions using statistical baselines, and order inventory based on past sales data.

At lunch you go to the No. 1 restaurant based on a survey of 500 people. You eat food that was priced based on marketing data. You go to your doctor's appointment where they take your blood pressure, temperature, weight, and do a blood test; after all the information is collected, you get a report showing your numbers and how you compare

to the statistical norms.

You head home in your car that's been serviced by a computer running statistical diagnostics. When you get home, you turn on the news and hear the latest crime statistics, see how the stock market performed, and discover how many people visited the zoo last week.

At night, you brush your teeth with toothpaste that's been statistically proven to fight cavities, read a few pages of your *New York Times* Best-Seller (based on statistical sales estimates), and go to sleep — only to start it all over again the next morning. But how can you be sure that all those statistics you encounter and depend on each day are correct? In Chapter 2, I discuss in more depth a few examples of how statistics is involved in our lives and workplaces, what its impact is, and how you can raise your awareness of it.



Some statistics are vague, inappropriate, or just plain wrong. You need to become more aware of the statistics you encounter each day and train your mind to stop and say “wait a minute!”, sift through the information, ask questions, and raise red flags when something’s not quite right. In Chapter 3, you see ways in which you can be misled by bad statistics and develop skills to think critically and identify problems before automatically believing results.

Like any other field, statistics has its own set of jargon, and I outline and explain some of the most commonly used statistical terms in Chapter 4. Knowing the language increases your ability to understand and communicate statistics at a higher level without being intimidated. It raises your credibility when you use precise terms to describe what’s wrong with a statistical result (and why). And your presentations involving statistical tables, graphs, charts, and analyses will be informational and effective. (Heck, if nothing else, you need the jargon because I use it throughout this book; don’t worry though, I always review it.)

In the next sections, you see how statistics is involved in each phase of the scientific method.

Designing Appropriate Studies

Everyone’s asking questions, from drug companies to biologists; from marketing analysts to the U.S. government. And ultimately, everyone will use statistics to help them answer their questions. In particular, many medical and psychological studies are done because someone wants to know the answer to a question. For example,

- ✓ Will this vaccine be effective in preventing the flu?

- ✓ What do Americans think about the state of the economy?
- ✓ Does an increase in the use of social networking Web sites cause depression in teenagers?

The first step after a research question has been formed is to design an effective study to collect data that will help answer that question. This step amounts to figuring out what process you'll use to get the data you need. In this section, I give an overview of the two major types of studies — surveys and experiments — and explore why it's so important to evaluate how a study was designed before you believe the results.

Surveys

An *observational study* is one in which data is collected on individuals in a way that doesn't affect them. The most common observational study is the survey. *Surveys* are questionnaires that are presented to individuals who have been selected from a population of interest. Surveys take many different forms: paper surveys sent through the mail, questionnaires on Web sites, call-in polls conducted by TV networks, phone surveys, and so on.



If conducted properly, surveys can be very useful tools for getting information. However, if not conducted properly, surveys can result in bogus information. Some problems include improper wording of questions, which can be misleading, lack of response by people who were selected to participate, or failure to include an entire group of the population. These potential problems mean a survey has to be well thought out before it's given.



Many researchers spend a great deal of time and money to do good surveys, and you'll know (by the criteria I discuss in Chapter 16) that you can trust them. However, as you are besieged with so many different types of surveys found in the media, in the workplace, and in many of your classes, you need to be able to quickly examine and critique how a survey was designed and conducted and be able to point out specific problems in a well-informed way. The tools you need for sorting through surveys are found in Chapter 16.

Experiments

An *experiment* imposes one or more treatments on the participants in such a way that clear comparisons can be made. After the treatments are applied, the responses are recorded. For example, to study the effect of drug dosage on blood pressure, one group may take 10 mg of the drug, and another group may take 20 mg. Typically, a control group

is also involved, in which subjects each receive a fake treatment (a sugar pill, for example), or a standard, nonexperimental treatment (like the existing drugs given to AIDS patients.)



Good and credible experiments are designed to minimize bias, collect lots of good data, and make appropriate comparisons (treatment group versus control group). Some potential problems that occur with experiments include researchers and/or subjects who know which treatment they got, factors not controlled for in the study that affect the outcome (such as weight of the subject when studying drug dosage), or lack of a control group (leaving no baseline to compare the results with).

But when designed correctly, an experiment can help a researcher establish a cause-and-effect relationship if the difference in responses between the treatment group and the control group is statistically significant (unlikely to have occurred just by chance).



Experiments are credited with helping to create and test drugs, determining best practices for making and preparing foods, and evaluating whether a new treatment can cure a disease, or at least reduce its impact. Our quality of life has certainly been improved through the use of well-designed experiments. However, not all experiments are well-designed, and your ability to determine which results are credible and which results are incredible (pun intended) is critical, especially when the findings are very important to you. All the info you need to know about experiments and how to evaluate them is found in Chapter 17.

Collecting Quality Data

After a study has been designed, be it a survey or an experiment, the individuals who will participate have to be selected, and a process must be in place to collect the data. This phase of the process is critical to producing credible data in the end, and this section hits the highlights.

Selecting a good sample



Statisticians have a saying, “Garbage in equals garbage out.” If you select your *subjects* (the individuals who will participate in your study) in a way that is *biased* — that is, favoring certain individuals or groups of individuals — then your results will also be biased. It’s that simple.

Suppose Bob wants to know the opinions of people in your city regarding a proposed

casino. Bob goes to the mall with his clipboard and asks people who walk by to give their opinions. What's wrong with that? Well, Bob is only going to get the opinions of a) people who shop at that mall; b) on that particular day; c) at that particular time; d) and who take the time to respond.

Those circumstances are too restrictive — those folks don't represent a cross section of the city. Similarly, Bob could put up a Web site survey and ask people to use it to vote. However, only people who know about the site, have Internet access, and want to respond will give him data, and typically only those with strong opinions will go to such trouble. In the end, all Bob has is a bunch of biased data on individuals that don't represent the city at all.



To minimize bias in a survey, the key word is *random*. You need to select your sample of individuals *randomly* — that is, with some type of “draw names out of a hat” process. Scientists use a variety of methods to select individuals at random, and you see how they do it in Chapter 16.

Note that in designing an experiment, collecting a random sample of people and asking them to participate often isn't ethical because experiments impose a treatment on the subjects. What you do is send out requests for volunteers to come to you. Then you make sure the volunteers you select from the group represent the population of interest and that the data is well collected on those individuals so the results can be projected to a larger group. You see how that's done in Chapter 17.

After going through Chapters 16 and 17, you'll know how to dig down and analyze others' methods for selecting samples and even be able to design a plan you can use to select a sample. In the end, you'll know when to say “Garbage in equals garbage out.”

Avoiding bias in your data

Bias is the systematic favoritism of certain individuals or certain responses. Bias is the nemesis of statisticians, and they do everything they can to minimize it. Want an example of bias? Say you're conducting a phone survey on job satisfaction of Americans; if you call people at home during the day between 9 a.m. and 5 p.m., you miss out on everyone who works during the day. Maybe day workers are more satisfied than night workers.

You have to watch for bias when collecting survey data. For instance: Some surveys are too long — what if someone stops answering questions halfway through? Or what if they give you misinformation and tell you they make \$100,000 a year instead of \$45,000? What if they give you answers that aren't on your list of possible answers? A host of problems can occur when collecting survey data, and you need to be able to pinpoint those problems.



Experiments are sometimes even more challenging when it comes to bias and collecting data. Suppose you want to test blood pressure; what if the instrument you're using breaks during the experiment? What if someone quits the experiment halfway through? What if something happens during the experiment to distract the subjects or the researchers? Or they can't find a vein when they have to do a blood test exactly one hour after a dose of a drug is given? These problems are just some examples of what can go wrong in data collection for experiments, and you have to be ready to look for and find these problems.

After you go through Chapter 16 (on samples and surveys) and Chapter 17 (on experiments), you'll be able to select samples and collect data in an unbiased way, being sensitive to little things that can really influence the results. And you'll have the ability to evaluate the credibility of statistical results and to be heard, because you'll know what you're talking about.

Creating Effective Summaries

After good data have been collected, the next step is to summarize them to get a handle on the big picture. Statisticians describe data in two major ways: with numbers (called *descriptive statistics*) and with pictures (that is, charts and graphs).

Descriptive statistics



Descriptive statistics are numbers that describe a data set in terms of its important features:

- ✓ If the data are *categorical* (where individuals are placed into groups, such as gender or political affiliation), they are typically summarized using the number of individuals in each group (called the *frequency*) or the percentage of individuals in each group (called the *relative frequency*).
- ✓ *Numerical data* represent measurements or counts, where the actual numbers have meaning (such as height and weight). With numerical data, more features can be summarized besides the number or percentage in each group. Some of these features include
 - Measures of center (in other words, where is the “middle” of the data?)
 - Measures of spread (how diverse or how concentrated are the data around the center?)
 - If appropriate, numbers that measure the relationship between two variables



Some descriptive statistics are more appropriate than others in certain situations; for example, the average isn't always the best measure of the center of a data set; the median is often a better choice. And the standard deviation isn't the only measure of variability on the block; the interquartile range has excellent qualities too. You need to be able to discern, interpret, and evaluate the types of descriptive statistics being presented to you on a daily basis and to know when a more appropriate statistic is in order.

The descriptive statistics you see most often are calculated, interpreted, compared, and evaluated in Chapter 5. These commonly used descriptive statistics include frequencies and relative frequencies (counts and percents) for categorical data; and the mean, median, standard deviation, percentiles, and their combinations for numerical data.

Charts and graphs

Data is summarized in a visual way using charts and/or graphs. These are displays that are organized to give you a big picture of the data in a flash and/or to zoom in on a particular result that was found. In this world of quick information and mini-sound bites, graphs and charts are commonplace. Most graphs and charts make their points clearly, effectively, and fairly; however, they can leave room for too much poetic license, and as a result, can expose you to a high number of misleading and incorrect graphs and charts.



In Chapters 6 and 7, I cover the major types of graphs and charts used to summarize both categorical and numerical data (see the preceding section for more about these types of data). You see how to make them, what their purposes are, and how to interpret the results. I also show you lots of ways that graphs and charts can be made to be misleading and how you can quickly spot the problems. It's a matter of being able to say "Wait a minute here! That's not right!" and knowing why not. Here are some highlights:

- ✓ Some of the basic graphs used for categorical data include pie charts and bar graphs, which break down variables, such as gender or which applications are used on teens' cellphones. A bar graph, for example, may display opinions on an issue using five bars labeled in order from "Strongly Disagree" up through "Strongly Agree." Chapter 6 gives you all the important info on making, interpreting, and, most importantly, evaluating these charts and graphs for fairness. You may be surprised to see how much can go wrong with a simple bar chart.
- ✓ For numerical data such as height, weight, time, or amount, a different type of

graph is needed. Graphs called histograms and boxplots are used to summarize numerical data, and they can be very informative, providing excellent on-the-spot information about a data set. But of course they also can be misleading, either by chance or even by design. (See Chapter 7 for the scoop.)



You're going to run across charts and graphs every day — you can open a newspaper and probably find several graphs without even looking hard. Having a statistician's magnifying glass to help you interpret the information is critical so that you can spot misleading graphs before you draw the wrong conclusions and possibly act on them. All the tools you need are ready for you in Chapter 6 (for categorical data) and Chapter 7 (for numerical data).

Determining Distributions

A *variable* is a characteristic that's being counted, measured, or categorized. Examples include gender, age, height, weight, or number of pets you own. A *distribution* is a listing of the possible values of a variable (or intervals of values), and how often (or at what density) they occur. For example, the distribution of gender at birth in the United States has been estimated at 52.4% male and 47.6% female.



Different types of distributions exist for different variables. The following three distributions are the most commonly occurring distributions in an introductory statistics course, and they have many applications in the real world:

- ✓ If a variable is counting the number of successes in a certain number of trials (such as the number of people who got well by taking a certain drug), it has a *binomial* distribution.
- ✓ If the variable takes on values that occur according to a “bell-shaped curve,” such as national achievement test scores, then that variable has a *normal* distribution.
- ✓ If the variable is based on sample averages and you have limited data, such as in a test of only ten subjects to see if a weight-loss program works, the *t*-distribution may be in order.

When it comes to distributions, you need to know how to decide which distribution a particular variable has, how to find probabilities for it, and how to figure out what the long-term average and standard deviation of the outcomes would be. To get you squared away on these issues, I've got three chapters for you, one dedicated to each distribution: Chapter 8 is all about the binomial, Chapter 9 handles the normal, and Chapter 10 focuses on the *t*-distribution.



For those of you taking an introductory statistics course (or any statistics course, for that matter), you know that one of the most difficult topics to understand is sampling distributions and the Central Limit Theorem (these two things go hand in hand). Chapter 11 walks you through these topics step by step so you understand what a sampling distribution is, what it's used for, and how it provides the foundation for data analyses like hypothesis tests and confidence intervals (see the next section for more about analyzing data). When you understand the Central Limit Theorem, it actually helps you solve difficult problems more easily, and all the keys to this information are there for you in Chapter 11.

Performing Proper Analyses

After the data have been collected and described using numbers and pictures, then comes the fun part: navigating through that black box called the *statistical analysis*. If the study has been designed properly, the original questions can be answered using the appropriate analysis — the operative word here being *appropriate*.



Many types of analyses exist, and choosing the right analysis for the right situation is critical, as is interpreting results properly, being knowledgeable of the limitations, and being able to evaluate others' choice of analyses and the conclusions they make with them.

In this book, you get all the information and tools you need to analyze data using the most common methods in introductory statistics: confidence intervals, hypothesis tests, correlation and regression, and the analysis of two-way tables. This section gives you a basic overview of those methods.

Margin of error and confidence intervals

You often see statistics that try to estimate numbers pertaining to an entire population; in fact, you see them almost every day in the form of survey results. The media tells you what the average gas price is in the U.S., how Americans feel about the job the president is doing, or how many hours people spend on the Internet each week.

But no one can give you a single-number result and claim it's an accurate estimate of the entire population unless he collected data on every single member of the population. For example, you may hear that 60 percent of the American people support the president's approach to healthcare, but you know they didn't ask you, so how could they have asked everybody? And since they didn't ask everybody, you know that a one-number answer isn't going to cut it.

What's really happening is that data is collected on a sample from the population (for example, the Gallup Organization calls 2,500 people at random), the results from that sample are analyzed, and conclusions are made regarding the entire population (for example, all Americans) based on those sample results.



The bottom line is, sample results vary from sample to sample, and this amount of variability needs to be reported (but it often isn't). The statistic used to measure and report the level of precision in someone's sample results is called the *margin of error*. In this context, the word *error* doesn't mean a mistake was made; it just means that because you didn't sample the entire population, a gap will exist between your results and the actual value you are trying to estimate for the population.

For example, someone finds that 60% of the 1,200 people surveyed support the president's approach to healthcare and reports the results with a margin of error of plus or minus 2%. This final result, in which you present your findings as a range of likely values between 58% and 62%, is called a *confidence interval*.



Everyone is exposed to results including a margin of error and confidence intervals, and with today's data explosion, many people are also using them in the workplace. Be sure you know what factors affect margin of error (like sample size) and what the makings of a good confidence interval are and how to spot them. You should also be able to find your own confidence intervals when you need to.

In Chapter 12, you find out everything you need to know about the margin of error: All the components of it, what it does and doesn't measure, and how to calculate it for a number of situations. Chapter 13 takes you step by step through the formulas, calculations, and interpretations of confidence intervals for a population mean, population proportion, and the difference between two means and proportions.

Hypothesis tests

One main staple of research studies is called hypothesis testing. A *hypothesis test* is a technique for using data to validate or invalidate a claim about a population. For example, a politician may claim that 80% of the people in her state agree with her — is that really true? Or, a company may claim that they deliver pizzas in 30 minutes or less; is that really true? Medical researchers use hypothesis tests all the time to test whether or not a certain drug is effective, to compare a new drug to an existing drug in terms of its side effects, or to see which weight-loss program is most effective with a certain group of people.



The elements about a population that are most often tested are

- ✓ The population mean (Is the average delivery time of 30 minutes really true?)
- ✓ The population proportion (Is it true that 80% of the voters support this candidate, or is it less than that?)
- ✓ The difference in two population means or proportions (Is it true that the average weight loss on this new program is 10 pounds more than the most popular program? Or, is it true that this drug decreases blood pressure by 10% more than the current drug?)



Hypothesis tests are used in a host of areas that affect your everyday life, such as medical studies, advertisements, polling data, and virtually anywhere that comparisons are made based on averages or proportions. And in the workplace, hypothesis tests are used heavily in areas like marketing, where you want to determine whether a certain type of ad is effective or whether a certain group of individuals buys more or less of your product now compared to last year.

Often you only hear the conclusions of hypothesis tests (for example, this drug is significantly more effective and has fewer side effects than the drug you are using now); but you don't see the methods used to come to these conclusions. Chapter 14 goes through all the details and underpinnings of hypothesis tests so you can conduct and critique them with confidence. Chapter 15 cuts right to the chase of providing step-by-step instructions for setting up and carrying out hypothesis tests for a host of specific situations (one population mean, one population proportion, the difference of two population means, and so on).

After reading Chapters 14 and 15, you'll be much more empowered when you need to know things like which group you should be marketing a product to; which brand of tires will last the longest; whether a certain weight-loss program is effective; and bigger questions like which surgical procedure you should opt for.

Correlation, regression, and two-way tables

One of the most common goals of research is to find links between variables. For example,

- ✓ Which lifestyle behaviors increase or decrease the risk of cancer?
- ✓ What side effects are associated with this new drug?
- ✓ Can I lower my cholesterol by taking this new herbal supplement?
- ✓ Does spending a large amount of time on the Internet cause a person to gain weight?

Finding links between variables is what helps the medical world design better drugs and treatments, provides marketers with info on who is more likely to buy their products, and gives politicians information on which to build arguments for and against certain policies.



In the mega-business of looking for relationships between variables, you find an incredible number of statistical results — but can you tell what's correct and what's not? Many important decisions are made based on these studies, and it's important to know what standards need to be met in order to deem the results credible, especially when a cause-and-effect relationship is being reported.

Chapter 18 breaks down all the details and nuances of plotting data from two numerical variables (such as dosage level and blood pressure), finding and interpreting *correlation* (the strength and direction of the linear relationship between x and y), finding the equation of a line that best fits the data (and when doing so is appropriate), and how to use these results to make predictions for one variable based on another (called *regression*). You also gain tools for investigating when a line fits the data well and when it doesn't, and what conclusions you can make (and shouldn't make) in the situations where a line does fit.

I cover methods used to look for and describe links between two categorical variables (such as the number of doses taken per day and the presence or absence of nausea) in detail in Chapter 19. I also provide info on collecting and organizing data into *two-way tables* (where the possible values of one variable make up the rows and the possible values for the other variable make up the columns), interpreting the results, analyzing the data from two-way tables to look for relationships, and checking for independence. And, as I do throughout this book, I give you strategies for critically examining results of these kinds of analyses for credibility.

Drawing Credible Conclusions



To perform statistical analyses, researchers use statistical software that depends on formulas. But formulas don't know whether they are being used properly, and they don't warn you when your results are incorrect. At the end of the day, computers can't tell you what the results mean; you have to figure it out.

Throughout this book you see what kinds of conclusions you can and can't make after the analysis has been done. The following sections provide an introduction to drawing appropriate conclusions.

Reeling in overstated results

Some of the most common mistakes made in conclusions are overstating the results or generalizing the results to a larger group than was actually represented by the study. For example, a professor wants to know which Super Bowl commercials viewers liked best. She gathers 100 students from her class on Super Bowl Sunday and asks them to rate each commercial as it is shown. A top-five list is formed, and she concludes that all Super Bowl viewers liked those five commercials the best. But she really only knows which ones *her students* liked best — she didn't study any other groups, so she can't draw conclusions about all viewers.

Questioning claims of cause and effect

One situation in which conclusions cross the line is when researchers find that two variables are related (through an analysis such as regression; see the earlier section “Correlation, regression, and two-way tables” for more info) and then automatically leap to the conclusion that those two variables have a cause-and-effect relationship.

For example, suppose a researcher conducted a health survey and found that people who took vitamin C every day reported having fewer colds than people who didn't take vitamin C every day. Upon finding these results, she wrote a paper and gave a press release saying vitamin C prevents colds, using this data as evidence.

Now, while it may be true that vitamin C does prevent colds, this researcher's study can't claim that. Her study was observational, which means she didn't control for any other factors that could be related to both vitamin C and colds. For example, people who take vitamin C every day may be more health conscious overall, washing their hands more often, exercising more, and eating better foods; all these behaviors may be helpful in reducing colds.



Until you do a controlled experiment, you can't make a cause-and-effect conclusion based on relationships you find. (I discuss experiments in more detail earlier in this chapter.)

Becoming a Sleuth, Not a Skeptic

Statistics is about much more than numbers. To really “get” statistics, you need to understand how to make appropriate conclusions from studying data and be savvy enough to not believe everything you hear or read until you find out how the information came about, what was done with it, and how the conclusions were drawn. That's something I discuss throughout the book, but I really zoom in on it in Chapter 20, which gives you ten ways to be a statistically savvy sleuth by recognizing common mistakes made by researchers and the media.



For you students out there, Chapter 21 brings good statistical practice into the exam setting and gives you tips on increasing your scores. Much of my advice is based on understanding the big picture as well as the details of tackling statistical problems and coming out a winner on the other side.



Becoming skeptical or cynical about statistics is very easy, especially after finding out what's going on behind the scenes; don't let that happen to you. You can find lots of good information out there that can affect your life in a positive way. Find a good channel for your skepticism by setting two personal goals:

- To become a well-informed consumer of the statistical information you see every day
- To establish job security by being the statistics "go-to" person who knows when and how to help others and when to find a statistician

Through reading and using the information in this book, you'll be confident in knowing you can make good decisions about statistical results. You'll conduct your own statistical studies in a credible way. And you'll be ready to tackle your next office project, critically evaluate that annoying political ad, or ace your next exam!

Chapter 2

The Statistics of Everyday Life

In This Chapter

- ▶ Raising questions about statistics you see in everyday life
 - ▶ Encountering statistics in the workplace
-

Today's society is completely taken over by numbers. Numbers are everywhere you look, from billboards showing the on-time statistics for a particular airline, to sports shows discussing the Las Vegas odds for upcoming football games. The evening news is filled with stories focusing on crime rates, the expected life span of junk-food junkies, and the president's approval rating. On a normal day, you can run into 5, 10, or even 20 different statistics (with many more on election night). Just by reading a Sunday newspaper all the way through, you come across literally hundreds of statistics in reports, advertisements, and articles covering everything from soup (how much does an average person consume per year?) to nuts (almonds are known to have positive health effects — what about other types of nuts?).

In this chapter I discuss the statistics that often appear in your life and work and talk about how statistics are presented to the general public. After reading this chapter, you'll realize just how often the media hits you with numbers and how important it is to be able to unravel the meaning of those numbers. Like it or not, statistics are a big part of your life. So, if you can't beat 'em, join 'em. And if you don't want to join 'em, at least try to understand 'em.

Statistics and the Media: More Questions than Answers?

Open a newspaper and start looking for examples of articles and stories involving numbers. It doesn't take long before numbers begin to pile up. Readers are inundated with results of studies, announcements of breakthroughs, statistical reports, forecasts, projections, charts, graphs, and summaries. The extent to which statistics occur in the media is mind-boggling. You may not even be aware of how many times you're hit with numbers nowadays.

This section looks at just a few examples from one Sunday paper's worth of news that I read the other day. When you see how frequently statistics are reported in the news without providing all the information you need, you may find yourself getting nervous,

wondering what you can and can't believe anymore. Relax! That's what this book is for — to help you sort out the good information from the bad (the chapters in Part II give you a great start on that).

Probing popcorn problems

The first article I came across that dealt with numbers was "Popcorn plant faces health probe," with the subheading: "Sick workers say flavoring chemicals caused lung problems." The article describes how the Centers for Disease Control (CDC) expressed concern about a possible link between exposure to chemicals in microwave popcorn flavorings and some cases of fixed obstructive lung disease. Eight people from one popcorn factory alone contracted this lung disease, and four of them were awaiting lung transplants.

According to the article, similar cases were reported at other popcorn factories. Now, you may be wondering, what about the folks who eat microwave popcorn? According to the article, the CDC finds "no reason to believe that people who eat microwave popcorn have anything to fear." (Stay tuned.) The next step is to evaluate employees more in-depth, including conducting surveys to determine health and possible exposures to the said chemicals, checks of lung capacity, and detailed air samples. The question here is: How many cases of this lung disease constitute a real pattern, compared to mere chance or a statistical anomaly? (You find out more about this in Chapter 14.)

Venturing into viruses

The second article discussed a recent cyber attack: A wormlike virus made its way through the Internet, slowing down Web browsing and e-mail delivery across the world. How many computers were affected? The experts quoted in the article said that 39,000 computers were infected, and they in turn affected hundreds of thousands of other systems.

Questions: How did the experts get that number? Did they check each computer out there to see whether it was affected? The fact that the article was written less than 24 hours after the attack suggests the number is a guess. Then why say 39,000 and not 40,000 — to make it seem less like a guess? To find out more on how to guesstimate with confidence (and how to evaluate someone else's numbers), see Chapter 13.

Comprehending crashes

Next in the paper was an alert about the soaring number of motorcycle fatalities. Experts said that the *fatality rate* — the number of fatalities per 100,000 registered vehicles — for motorcyclists has been steadily increasing, as reported by the National Highway Traffic Safety Administration (NHTSA). In the article, many possible causes for the increased

motorcycle death rate are discussed, including age, gender, size of engine, whether the driver had a license, alcohol use, and state helmet laws (or lack thereof). The report is very comprehensive, showing various tables and graphs with the following titles:

- ✓ Motorcyclists killed and injured, and fatality and injury rates by year, per number of registered vehicles, and per millions of vehicle miles traveled
- ✓ Motorcycle rider fatalities by state, helmet use, and blood alcohol content
- ✓ Occupant fatality rates by vehicle type (motorcycles, passenger cars, light trucks), per 10,000 registered vehicles and per 100 million vehicle miles traveled
- ✓ Motorcyclist fatalities by age group
- ✓ Motorcyclist fatalities by engine size (displacement)
- ✓ Previous driving records of drivers involved in fatal traffic crashes by type of vehicle (including previous crashes, DUI convictions, speeding convictions, and license suspensions and revocations)
- ✓ Percentage of alcohol-impaired motorcycle riders killed in traffic crashes by time of day, for single-vehicle, multiple-vehicle, and total crashes

This article is very informative and provides a wealth of detailed information regarding motorcycle fatalities and injuries in the U.S. However, the onslaught of so many tables, graphs, rates, numbers, and conclusions can be overwhelming and confusing and allow you to miss the big picture. With a little practice, and help from Part II, you'll be better able to sort out graphs, tables, and charts and all the statistics that go along with them. For example, some important statistical issues come up when you see rates versus counts (such as death rates versus number of deaths). As I address in Chapter 3, counts can give you misleading information if they're used when rates would be more appropriate.

Mulling malpractice

Further along in the newspaper was a report about a recent medical malpractice insurance study: Malpractice cases affect people in terms of the fees doctors charge and the ability to get the healthcare they need. The article indicates that 1 in 5 Georgia doctors have stopped doing risky procedures (such as delivering babies) because of the ever-increasing malpractice insurance rates in the state. This is described as a "national epidemic" and a "health crisis" around the country. Some brief details of the study are included, and the article states that of the 2,200 Georgia doctors surveyed, 2,800 of them — which they say represents about 18% of those sampled — were expected to stop providing high-risk procedures.

Wait a minute! That can't be right. Out of 2,200 doctors, 2,800 don't perform the

procedures, and that is supposed to represent 18%? That's impossible! You can't have a bigger number on the top of a fraction, and still have the fraction be under 100%, right? This is one of many examples of errors in media reporting of statistics. So what's the real percentage? There's no way to tell from the article. Chapter 5 nails down the particulars of calculating statistics so that you can know what to look for and immediately tell when something's not right.

Belaboring the loss of land

In the same Sunday paper was an article about the extent of land development and speculation across the United States. Knowing how many homes are likely to be built in your neck of the woods is an important issue to get a handle on. Statistics are given regarding the number of acres of farmland being lost to development each year. To further illustrate how much land is being lost, the area is also listed in terms of football fields. In this particular example, experts said that the mid-Ohio area is losing 150,000 acres per year, which is 234 square miles, or 115,385 football fields (including end zones). How do people come up with these numbers, and how accurate are they? And does it help to visualize land loss in terms of the corresponding number of football fields? I discuss the accuracy of data collected in more detail in Chapter 16.

Scrutinizing schools

The next topic in the paper was school proficiency — specifically, whether extra school sessions help students perform better. The article states that 81.3% of students in this particular district who attended extra sessions passed the writing proficiency test, whereas only 71.7% of those who didn't participate in the extra school sessions passed it. But is this enough of a difference to account for the \$386,000 price tag per year? And what's happening in these sessions to cause an improvement? Are students in these sessions spending more time just preparing for those exams rather than learning more about writing in general? And here's the big question: Were the participants in the extra sessions student volunteers who may be more motivated than the average student to try to improve their test scores? The article doesn't say.

Studying surveys of all shapes and sizes

Surveys and polls are among the most visible mechanisms used by today's media to grab your attention. It seems that everyone wants to do a survey, including market managers, insurance companies, TV stations, community groups, and even students in high school classes. Here are just a few examples of survey results that are part of today's news:

With the aging of the American workforce, companies are planning for their future leadership. (How do they know that the American workforce is aging, and if it is, by how much is it aging?) A recent survey shows that nearly 67% of human-resources managers polled said that planning for succession had become more important in the past five years than it had been in the past. The survey also says that

88% of the 210 respondents said they usually or often fill senior positions with internal candidates. But how many managers did not respond, and is 210 respondents really enough people to warrant a story on the front page of the business section? Believe it or not, when you start looking for them, you'll find numerous examples in the news of surveys based on far fewer participants than 210. (To be fair, however, 210 can actually be a good number of subjects in some situations. The issues of what sample size is large enough and what percentage of respondents is big enough are addressed in full detail in Chapter 16.)

Some surveys are based on current interests and trends. For example, a recent Harris-Interactive survey found that nearly half (47%) of U.S. teens say their social lives would end or be worsened without their cellphones, and 57% go as far as to say that their cellphones are the key to their social life. The study also found that 42% of teens say that they can text while blindfolded (how do you really test this?). Keep in perspective, though, that the study did not tell you what percentage of teens actually have cellphones or what demographic characteristics those teens have compared to teens who do *not* have cellphones. And remember that data collected on topics like this aren't always accurate, because the individuals who are surveyed may tend to give biased answers (who wouldn't want to say they can text blindfolded?). For more information on how to interpret and evaluate the results of surveys, see Chapter 16.

Studies like this appear all the time, and the only way to know what to believe is to understand what questions to ask and to be able to critique the quality of the study. That's all part of statistics! The good news is, with a few clarifying questions, you can quickly critique statistical studies and their results. Chapter 17 helps you do just that.

Studying sports

The sports section is probably the most numerically jampacked section of the newspaper. Beginning with game scores, the win/loss percentages for each team, and the relative standing for each team, the specialized statistics reported in the sports world are so deep they require wading boots to get through. For example, basketball statistics are broken down by team, by quarter, and by player. For each player, you get minutes played, field goals, free throws, rebounds, assists, personal fouls, turnovers, blocks, steals, and total points.

Who needs to know this stuff, besides the players' mothers? Apparently many fans do. Statistics are something that sports fans can never get enough of and players often can't stand to hear about. Stats are the substance of water-cooler debates and the fuel for armchair quarterbacks around the world.

Fantasy sports have also made a huge impact on the sports money-making machine. Fantasy sports are games where participants act as owners to build their own teams from existing players in a professional league. The fantasy team owners then compete against each other. What is the competition based on? Statistical performance of the players and

teams involved, as measured by rules set up by a “league commissioner” and an established point system. According to the Fantasy Sports Trade Association, the number of people age 12 and up who are involved in fantasy sports is more than 30 million, and the amount of money spent is \$3–4 billion per year. (And even here you can ask how the numbers were calculated — the questions never end, do they?)

Banking on business news

The business section of the newspaper provides statistics about the stock market. In one week the market went down 455 points; is that decrease a lot or a little? You need to calculate a percentage to really get a handle on that.

The business section of my paper contained reports on the highest yields nationwide on every kind of certificate of deposit (CD) imaginable. (By the way, how do they know those yields are the highest?) I also found reports about rates on 30-year fixed loans, 15-year fixed loans, 1-year adjustable rate loans, new car loans, used car loans, home equity loans, and loans from your grandmother (well actually no, but if grandma read these statistics, she might increase her cushy rates).

Finally, I saw numerous ads for those beloved credit cards — ads listing the interest rates, the annual fees, and the number of days in the billing cycle. How do you compare all the information about investments, loans, and credit cards in order to make a good decision? What statistics are most important? The real question is: Are the numbers reported in the paper giving the whole story, or do you need to do more detective work to get at the truth? Chapters 16 and 17 help you start tearing apart these numbers and making decisions about them.

Touring the travel news

You can’t even escape the barrage of numbers by heading to the travel section. For example, there I found that the most frequently asked question coming in to the Transportation Security Administration’s response center (which receives about 2,000 telephone calls, 2,500 e-mail messages, and 200 letters per week on average — would you want to be the one counting all of those?) is, “Can I carry this on a plane?” This can refer to anything from an animal to a wedding dress to a giant tin of popcorn. (I wouldn’t recommend the tin of popcorn. You have to put it in the overhead compartment horizontally, and because things shift during flight, the cover will likely open; and when you go to claim your tin at the end of the flight, you and your seatmates will be showered. Yes, I saw it happen once.)

The number of reported responses in this case leads to an interesting statistical question: How many operators are needed at various times of the day to field those calls, e-mails, and letters coming in? Estimating the number of anticipated calls is your first step, and being wrong can cost you money (if you overestimate it) or a lot of bad PR (if

you underestimate it). These kinds of statistical challenges are tackled in Chapter 13.

Surveying sexual stats

In today's age of info-overkill, it's very easy to find out what the latest buzz is, including the latest research on people's sex lives. An article in my paper reported that married people have 6.9 more sexual encounters per year than people who have never been married. That's nice to know, I guess, but how did someone come up with this number? The article I'm looking at doesn't say (maybe some statistics are better left unsaid?).

If someone conducted a survey by calling people on the phone asking for a few minutes of their time to discuss their sex lives, who will be the most likely to want to talk about it? And what are they going to say in response to the question, "How many times a week do you have sex?" Are they going to report the honest truth, tell you to mind your own business, or exaggerate a little? Self-reported surveys can be a real source of bias and can lead to misleading statistics. But how would you recommend people go about finding out more about this very personal subject? Sometimes, research is more difficult than it seems. (Chapter 16 discusses biases that come up when collecting certain types of survey data.)

Breaking down weather reports

Weather reports provide another mass of statistics, with forecasts of the next day's high and low temperatures (how do they decide it'll be 16 degrees and not 15 degrees?) along with reports of the day's UV factor, pollen count, pollution standard index, and water quality and quantity. (How do they get these numbers — by taking samples? How many samples do they take, and where do they take them?) You can find out what the weather is right now anywhere in the world. You can get a forecast looking ahead three days, a week, a month, or even a year! Meteorologists collect and record tons and tons of data on the weather each day. Not only do these numbers help you decide whether to take your umbrella to work, but they also help weather researchers to better predict longer term forecasts and even global climate changes over time.

Even with all the information and technologies available to weather researchers, how accurate are weather reports these days? Given the number of times you get rained on when you were told it was going to be sunny, it seems they still have work to do on those forecasts. What the abundance of data really shows though, is that the number of variables affecting weather is almost overwhelming, not just to you, but for meteorologists, too.



Statistical computer models play an important role in making predictions about major weather-related events, such as hurricanes, earthquakes, and volcano

eruptions. Scientists still have some work to do before they can predict tornados before they begin to form or tell you exactly where and when a hurricane is going to hit land, but that's certainly their goal, and they continue to get better at it. For more on modeling and statistics, see Chapter 18.

Musing about movies

Moving on to the arts section, I saw several ads for current movies. Each movie ad contains quotes from certain movie critics: "Two thumbs up!" "The supreme adventure of our time," "Absolutely hilarious," or "One of the top ten films of the year!" Do you pay attention to the critics? How do you determine which movies to go to? Experts say that although the popularity of a movie may be affected by the critics' comments (good or bad) in the beginning of a film's run, word of mouth is the most important determinant of how well a film does in the long run.

Studies also show that the more dramatic a movie is, the more popcorn is sold. Yes, the entertainment business even keeps tabs on how much crunching you do at the movies. How do they collect all this information, and how does it impact the types of movies that are made? This, too, is part of statistics: designing and carrying out studies to help pinpoint an audience and find out what they like, and then using the information to help guide the making of the product. So the next time someone with a clipboard asks if you have a minute, you may want to stand up and be counted.

Highlighting horoscopes

Those horoscopes: You read them, but do you believe them? Should you? Can people predict what will happen more often than just by chance? Statisticians have a way of finding out, by using something they call a *hypothesis test* (see Chapter 14). So far they haven't found anyone who can read minds, but people still keep trying!

Using Statistics at Work

Now put down the Sunday newspaper and move on to the daily grind of the workplace. If you're working for an accounting firm, of course numbers are part of your daily life. But what about people like nurses, portrait studio photographers, store managers, newspaper reporters, office staff, or construction workers? Do numbers play a role in those jobs? You bet. This section gives you a few examples of how statistics creep into *every* workplace.



You don't have to go far to see how statistics weaves its way in and out of your life

and work. The secret is being able to determine what it all means and what you can believe, and to be able to make sound decisions based on the real story behind numbers so you can handle and become used to the statistics of everyday life.

Delivering babies — and information

Sue works as a nurse during the night shift in the labor and delivery unit at a university hospital. She takes care of several patients in a given evening, and she does her best to accommodate everyone. Her nursing manager has told her that each time she comes on shift she should identify herself to the patient, write her name on the whiteboard in the patient's room, and ask whether the patient has any questions. Why? Because a few days after each mother leaves with her baby, the hospital gives her a phone call asking about the quality of care, what was missed, what it could do to improve its service and quality of care, and what the staff could do to ensure that the hospital is chosen over other hospitals in town. For example, surveys show that patients who know the names of their nurses feel more comfortable, ask more questions, and have a more positive experience in the hospital than those who don't know the names of their nurses. Sue's salary raises depend on her ability to follow through with the needs of new mothers. No doubt the hospital has also done a lot of research to determine the factors involved in quality of patient care well beyond nurse-patient interactions. (See Chapter 17 for in-depth info concerning medical studies.)

Posing for pictures

Carol recently started working as a photographer for a department store portrait studio; one of her strengths is working with babies. Based on the number of photos purchased by customers over the years, this store has found that people buy more posed pictures than natural-looking ones. As a result, store managers encourage their photographers to take posed shots.

A mother comes in with her baby and has a special request: "Could you please not pose my baby too deliberately? I just like his pictures to look natural." If Carol says, "Can't do that, sorry. My raises are based on my ability to pose a child well," you can bet that the mother is going to fill out that survey on quality service after this session — and not just to get \$2.00 off her next sitting (if she ever comes back). Instead, Carol should show her boss the information in Chapter 16 about collecting data on customer satisfaction.

Poking through pizza data

Terry is a store manager at a local pizzeria that sells pizza by the slice. He is in charge of determining how many workers to have on staff at a given time, how many pizzas to make ahead of time to accommodate the demand, and how much cheese to order and grate, all

with minimal waste of wages and ingredients. Friday night at midnight, the place is dead. Terry has five workers left and has five large pans of pizza he could throw in the oven, making about 40 slices of pizza each. Should he send two of his workers home? Should he put more pizza in the oven or hold off?

The store owner has been tracking the demand for weeks now, so Terry knows that every Friday night things slow down between 10 and 12 p.m., but then the bar crowd starts pouring in around midnight and doesn't let up until the doors close at 2:30 a.m. So Terry keeps the workers on, puts in the pizzas in 30-minute intervals from midnight on, and is rewarded with a profitable night, with satisfied customers and with a happy boss. For more information on how to make good estimates using statistics, see Chapter 13.

Statistics in the office

D.J. is an administrative assistant for a computer company. How can statistics creep into her office workplace? Easy. Every office is filled with people who want to know answers to questions, and they want someone to "Crunch the numbers," to "Tell me what this means," to "Find out if anyone has any hard data on this," or to simply say, "Does this number make any sense?" They need to know everything from customer satisfaction figures to changes in inventory during the year; from the percentage of time employees spend on e-mail to the cost of supplies for the last three years. Every workplace is filled with statistics, and D.J.'s marketability and value as an employee could go up if she's the one the head honchos turn to for help. Every office needs a resident statistician — why not let it be you?

Chapter 3

Taking Control: So Many Numbers, So Little Time

In This Chapter

- ▶ Examining the extent of statistics abuse
 - ▶ Feeling the impact of statistics gone wrong
-

The sheer amount of statistics in daily life can leave you feeling overwhelmed and confused. This chapter gives you a tool to help you deal with statistics: skepticism! Not radical skepticism like “I can’t believe anything anymore,” but healthy skepticism like “Hmm, I wonder where that number came from?” and “I need to find out more information before I believe these results.” To develop healthy skepticism, you need to understand how the chain of statistical information works.

Statistics end up on your TV and in your newspaper as a result of a process. First, the researchers who study an issue generate results; this group is composed of pollsters, doctors, marketing researchers, government researchers, and other scientists. They are considered the *original sources* of the statistical information.

After they get their results, these researchers naturally want to tell people about it, so they typically either put out a press release or publish a journal article. Enter the journalists or reporters, who are considered the *media sources* of the information. Journalists hunt for interesting press releases and sort through journals, basically searching for the next headline. When reporters complete their stories, statistics are immediately sent out to the public through all forms of media. Now the information is ready to be taken in by the third group — the *consumers* of the information (you). You and other consumers of information are faced with the task of listening to and reading the information, sorting through it, and making decisions about it.

At any stage in the process of doing research, communicating results, or consuming information, errors can take place, either unintentionally or by design. The tools and strategies you find in this chapter give you the skills to be a good detective.

Detecting Errors, Exaggerations, and Just Plain Lies

Statistics can go wrong for many different reasons. First, a simple, honest error can occur. This can happen to anyone, right? Other times, the error is something other than a

simple, honest mistake. In the heat of the moment, because someone feels strongly about a cause and because the numbers don't quite bear out the point that the researcher wants to make, statistics get tweaked, or, more commonly, exaggerated, either in their values or how they're represented and discussed.

Another type of error is an *error of omission* — information that is missing that would have made a big difference in terms of getting a handle on the real story behind the numbers. That omission makes the issue of correctness difficult to address, because you're lacking information to go on.

You may even encounter situations in which the numbers have been completely fabricated and can't be repeated by anyone because they never happened. This section gives you tips to help you spot errors, exaggerations, and lies, along with some examples of each type of error that you, as an information consumer, may encounter.

Checking the math

The first thing you want to do when you come upon a statistic or the result of a statistical study is to ask, "Is this number correct?" Don't assume it is! You'd probably be surprised at the number of simple arithmetic errors that occur when statistics are collected, summarized, reported, or interpreted.



To spot arithmetic errors or omissions in statistics:

- ✓ **Check to be sure everything adds up.** In other words, do the percents in the pie chart add up to 100 (or close enough due to rounding)? Do the number of people in each category add up to the total number surveyed?
- ✓ **Double-check even the most basic calculations.**
- ✓ **Always look for a total so you can put the results into proper perspective.** Ignore results based on tiny sample sizes.
- ✓ **Examine whether the projections are reasonable.** For example, if three deaths due to a certain condition are said to happen per minute, that adds up to over 1.5 million such deaths in a year. Depending on what condition is being reported, this number may be unreasonable.

Uncovering misleading statistics

By far, the most common abuses of statistics are subtle, yet effective, exaggerations of the truth. Even when the math checks out, the underlying statistics themselves can be misleading if they exaggerate the facts. Misleading statistics are harder to pinpoint than

simple math errors, but they can have a huge impact on society, and, unfortunately, they occur all the time.

Breaking down statistical debates

Crime statistics are a great example of how statistics are used to show two sides of a story, only one of which is really correct. Crime is often discussed in political debates, with one candidate (usually the incumbent) arguing that crime has gone down during her tenure, and the challenger often arguing that crime has gone up (giving the challenger something to criticize the incumbent for). How can two candidates make such different conclusions based on the same data set? Turns out, depending on the way you measure crime, getting either result can be possible.

Table 3-1 shows the population of the United States for 1998 to 2008, along with the number of reported crimes and the crime *rates* (crimes per 100,000 people), calculated by taking the number of crimes divided by the population size and multiplying by 100,000.

Table 3-1 Number of Crimes, Estimated Population Size, and Crime Rates in the U.S.

Year	No. of Crimes	Population Size	Crime Rate per 100,000 People
1998	12,475,634	270,296,000	4,615.5
1999	11,634,378	272,690,813	4,266.5
2000	11,608,072	281,421,906	4,124.8
2001	11,876,669	285,317,559	4,162.6
2002	11,878,954	287,973,924	4,125.0
2003	11,826,538	290,690,788	4,068.4
2004	11,679,474	293,656,842	3,977.3
2005	11,565,499	296,507,061	3,900.6
2006	11,401,511	299,398,484	3,808.1
2007	11,251,828	301,621,157	3,730.5
2008	11,149,927	304,059,784	3,667.0

Source: U.S. Crime Victimization Survey

Now compare the number of crimes and the crime rates for 2001 and 2002 in Table 3-1. In column 2, you see that the *number of crimes* increased by 2,285 from 2001 to 2002 ($11,878,954 - 11,876,669$). This represents an increase of 0.019% (dividing the difference, 2,285, by the number of crimes in 2001, 11,876,669). Note the population size (column 3) also increased from 2001 to 2002, by 2,656,365 people ($287,973,924 - 285,317,559$), or 0.931% (dividing this difference by the population size in 2001). However, in column 4, you see the crime *rate* decreased from 2001 to 2002 from 4,162.6 (per 100,000 people) in 2001 to 4,125.0 (per 100,000) in 2002. How did the crime rate decrease? Although the number of crimes and the number of people both went up, the number of crimes increased at a slower rate than the increase in population size did (0.019% compare to

0.931%).

So how should the crime trend be reported? Did crime actually go up or down from 2001 to 2002? Based on the crime rate — which is a more accurate gauge — you can conclude that crime decreased during that year. But be watchful of the politician who wants to show that the incumbent didn't do his job; he will be tempted to look at the number of crimes and claim that crime went up, creating an artificial controversy and resulting in confusion (not to mention skepticism) on behalf of the voters. (Aren't election years fun?)



To create an even playing field when measuring how often an event occurs, you convert each number to a percent by dividing by the total to get what statisticians call a *rate*. Rates are usually better than count data because rates allow you to make fair comparisons when the totals are different.

Untwisting tornado statistics

Which state has the most tornados? It depends on how you look at it. If you just count the number of tornados in a given year (which is how I've seen the media report it most often), the top state is Texas. But think about it. Texas is the second biggest state (after Alaska). Yes, Texas is in that part of the U.S. called "Tornado Alley," and yes, it gets a lot of tornados, but it also has a huge surface area for those tornados to land and run.

A more fair comparison, and how meteorologists look at it, is to look at the number of tornados per 10,000 square miles. Using this statistic (depending on your source), Florida comes out on top, followed by Oklahoma, Indiana, Iowa, Kansas, Delaware, Louisiana, Mississippi, and Nebraska, and finally Texas weighs in at number 10. (Although I'm sure this is one statistic they are happy to rank low on; as opposed to their AP rankings in NCAA football.)

Other tornado statistics measured and reported include the state with the highest percentage of killer tornadoes as a percentage of all tornados (Tennessee); and the total length of tornado paths per 10,000 square miles (Mississippi). Note each of these statistics is reported appropriately as a *rate* (amount per unit).



Before believing statistics indicating "the highest XXX" or "the lowest XXX," take a look at how the variable is measured to see whether it's fair and whether there are other statistics that should be examined too to get the whole picture. Also make sure the units are appropriate for making fair comparisons.

Zeroing in on what the scale tells you

Charts and graphs are useful for making a quick and clear point about your data.

Unfortunately, many times the charts and graphs accompanying everyday statistics aren't done correctly and/or fairly. One of the most important elements to watch for is the way that the chart or graph is scaled. The *scale* of a graph is the quantity used to represent each tick mark on the axis of the graph. Do the tick marks increase by 1s, 10s, 20s, 100s, 1,000s, or what? The scale can make a big difference in terms of the way the graph or chart looks.

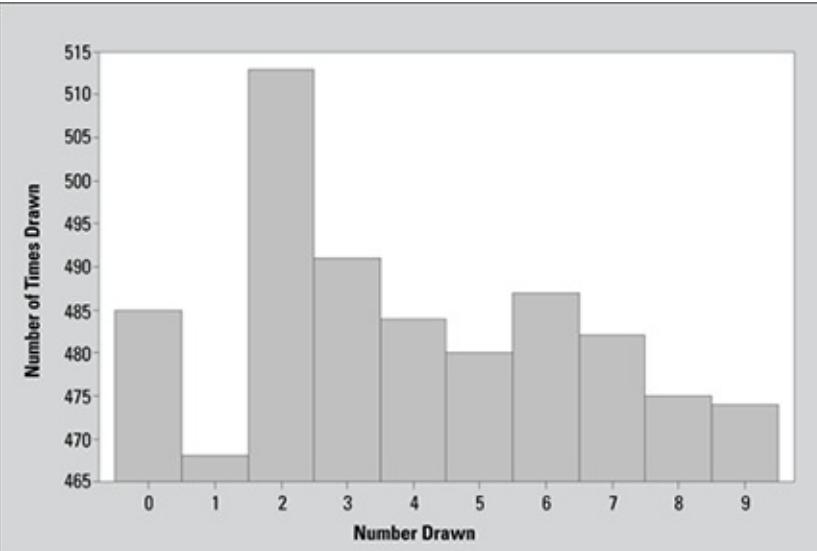
For example, the Kansas Lottery routinely shows its recent results from the Pick 3 Lottery. One of the statistics reported is the number of times each number (0 through 9) is drawn among the three winning numbers. Table 3-2 shows a chart of the number of times each number was drawn during 1,613 total Pick 3 games (4,839 single numbers drawn). It also reports the percentage of times that each number was drawn. Depending on how you choose to look at these results, you can again make the statistics appear to tell very different stories.

Table 3-2 **Numbers Drawn in the Pick 3 Lottery**

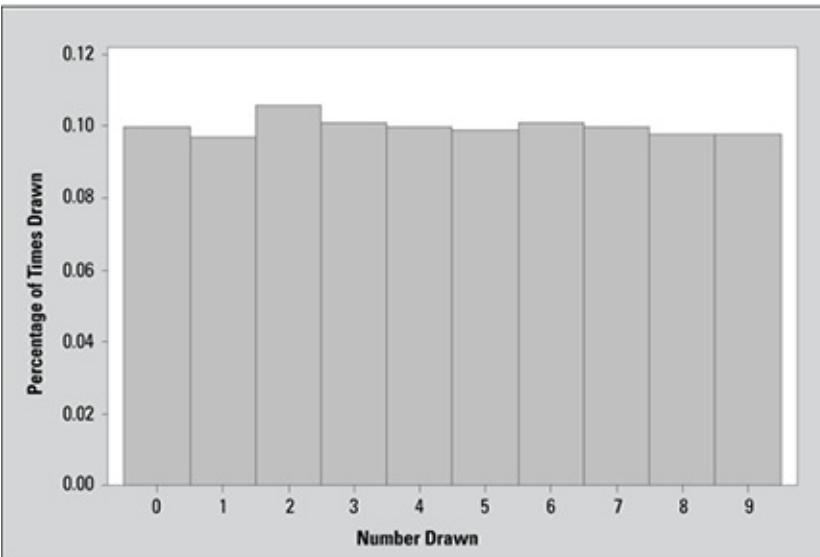
Number Drawn	No. of Times Drawn out of 4,839	Percentage of Times Drawn (No. of Times Drawn ÷ 4,839)
0	485	10.0%
1	468	9.7%
2	513	10.6%
3	491	10.1%
4	484	10.0%
5	480	9.9%
6	487	10.1%
7	482	10.0%
8	475	9.8%
9	474	9.8%

The way lotteries typically display results like those in Table 3-2 is shown in Figure 3-1a. Notice that in this chart, it seems that the number 1 doesn't get drawn nearly as often (only 468 times) as number 2 does (513 times). The difference in the height of these two bars appears to be very large, exaggerating the difference in the number of times these two numbers were drawn. However, to put this in perspective, the actual difference here is $513 - 468 = 45$ out of a total of 4,839 numbers drawn. In terms of percentages, the difference between the number of times the number 1 and the number 2 are drawn is $45 \div 4,839 = 0.009$, or only nine-tenths of one percent.

drawn; and b) percentage of times each number was drawn.



a



b

What makes this chart exaggerate the differences? Two issues come to mind. First, notice that the vertical axis, which shows the number of times (or frequency) that each number is drawn, goes up by 5s. So a difference of 5 out of a total of 4,839 numbers drawn appears significant. Stretching the scale so that differences appear larger than they really are is a common trick used to exaggerate results. Second, the chart starts counting at 465, not at 0. Only the top part of each bar is shown, which also exaggerates the results. In comparison, Figure 3-1b graphs the *percentage* of times each number was drawn. Normally the shape of a graph wouldn't change when going from counts to percentages; however, this chart uses a more realistic scale than the one in Figure 3-1a (going by 2% increments) and starts at 0, both of which make the differences appear as they really are — not much different at all. Boring, huh?

Maybe the lottery folks thought so too. In fact, maybe they use Figure 3-1a rather than Figure 3-1b because they want you to think that some "magic" is involved in the numbers, and you can't blame them; that's their business.



Looking at the scale of a graph or chart can really help you keep the reported

results in proper perspective. Stretching the scale out or starting the y-axis at the highest possible number makes differences appear larger; squeezing down the scale or starting the y-axis at a much lower value than needed makes differences appear smaller than they really are.

Checking your sources

When examining the results of any study, check the source of the information. The best results are often published in reputable journals that are well known by the experts in the field. For example, in the world of medical science, the *Journal of the American Medical Association* (JAMA), the *New England Journal of Medicine*, *The Lancet*, and the *British Medical Journal* are all reputable journals doctors use to publish results and read about new findings.



Consider the source and who financially supported the research. Many companies finance research and use it for advertising their products. Although that in itself isn't necessarily a bad thing, in some cases a conflict of interest on the part of researchers can lead to biased results. And if the results are very important to you, ask whether more than one study was conducted, and if so, ask to examine all the studies that were conducted, not just those whose results were published in journals or appeared in advertisements.

Counting on sample size

Sample size isn't everything, but it does count for a great deal in surveys and studies. If the study is designed and conducted correctly, and if the participants are selected randomly (that is, with no bias; see Chapter 16 for more on random samples), sample size is an important factor in determining the accuracy and repeatability of the results. (See Chapters 16 and 17 for more information on designing and carrying out studies.)

Many surveys are based on large numbers of participants, but that isn't always true for other types of research, such as carefully controlled experiments. Because of the high cost of some types of research in terms of time and money, some studies are based on a small number of participants or products. Researchers have to find the appropriate balance when determining sample size.



The most unreliable results are those based on *anecdotes*, stories that talk about a single incident in an attempt to sway opinion. Have you ever told someone not to buy a product because you had a bad experience with it? Remember that an anecdote (or story) is really a nonrandom sample whose size is only one.

Considering cause and effect

Headlines often simplify or skew the “real” information, especially when the stories involve statistics and the studies that generated the statistics.

A study conducted a few years back evaluated videotaped sessions of 1,265 patient appointments with 59 primary-care physicians and 6 surgeons in Colorado and Oregon. This study found that physicians who had not been sued for malpractice spent an average of 18 minutes with each patient, compared to 16 minutes for physicians who *had* been sued for malpractice. The study was reported by the media with the headline, “Bedside manner fends off malpractice suits.” However, this study seemed to say that if you are a doctor who gets sued, all you have to do is spend more time with your patients, and you’re off the hook. (Now when did bedside manner get characterized as time spent?)

Beyond that, are we supposed to believe that a doctor who has been sued needs only add a couple more minutes of time with each patient to avoid being sued in the future? Maybe what the doctor does during that time counts much more than how much time the doctor actually spends with each patient. You tackle the issues of cause-and-effect relationships between variables in Chapter 18.

Finding what you wanted to find

You may wonder how two political candidates can discuss the same topic and get two opposing conclusions, both based on “scientific surveys.” Even small differences in a survey can create big differences in results. (See Chapter 16 for the full scoop on surveys.)

One common source of skewed survey results comes from question wording. Here are three different questions that are trying to get at the same issue — public opinion regarding the line-item veto option available to the president:

- ✓ Should the line-item veto be available to the president to eliminate waste (yes/no/no opinion)?
- ✓ Does the line-item veto give the president too much individual power (yes/no/no opinion)?
- ✓ What is your opinion on the presidential line-item veto? Choose 1–5, with 1 = strongly opposed and 5 = strongly support.

The first two questions are misleading and will lead to biased results in opposite directions. The third version will draw results that are more accurate in terms of what people really think. However, not all surveys are written with the purpose of finding the truth; many are written to support a certain viewpoint.



Research shows that even small changes in wording affect survey outcomes,

leading to results that conflict when different surveys are compared. If you can tell from the wording of the question how they want you to respond to it, you know you're looking at a leading question; and leading questions lead to biased results. (See Chapter 16 for more on spotting problems with surveys.)

Looking for lies in all the right places

Every once in a while, you hear about someone who faked his data, or “fudged the numbers.” Probably the most commonly committed lie involving statistics and data is when people throw out data that don’t fit their hypothesis, don’t fit the pattern, or appear to be outliers. In cases when someone has clearly made an error (for example, someone’s age is recorded as 200), removing that erroneous data point or trying to correct the error makes sense. Eliminating data for any other reason is ethically wrong; yet it happens.

Regarding missing data from experiments, a commonly used phrase is “Among those who completed the study. . . .” What about those who didn’t complete the study, especially a medical one? Did they get tired of the side effects of the experimental drug and quit? If so, the loss of this person will create results that are biased toward positive outcomes.



Before believing the results of a study, check out how many people were chosen to participate, how many finished the study, and what happened to all the participants, not just the ones who experienced a positive result.

Surveys are not immune to problems from missing data, either. For example, it’s known by statisticians that the opinions of people who respond to a survey can be very different from the opinions of those who don’t. In general, the lower the percentage of people who respond to a survey (the response rate), the less credible the results will be. For more about surveys and missing data, see Chapter 16.

Feeling the Impact of Misleading Statistics

You make decisions every day based on statistics and statistical studies that you’ve heard about or seen, many times without even realizing it. Misleading statistics affect your life in small or large ways, depending on the type of statistics that cross your path and what you choose to do with the information you’re given. Here are some little everyday scenarios where statistics slip in:

- ✓ “Gee, I hope Rex doesn’t chew up my rugs again while I’m at work. I heard somewhere that dogs on Prozac deal better with separation anxiety. How did they

figure that out? And what would I tell my friends?"

- ✓ "I thought everyone was supposed to drink eight glasses of water a day, but now I hear that too much water could be bad for me; what should I believe?"
- ✓ "A study says people spend two hours a day at work checking and sending personal e-mails. How is that possible? No wonder my boss is paranoid."

You may run into other situations involving statistics that can have a larger impact on your life, and you need to be able to sort it all out. Here are some examples:

- ✓ A group lobbying for a new skateboard park tells you 80% of the people surveyed agree that taxes should be raised to pay for it, so you should too. Will you feel the pressure to say yes?
- ✓ The radio news at the top of the hour says cellphones cause brain tumors. Your spouse uses his cellphone all the time. Should you panic and throw away all cellphones in your house?
- ✓ You see an advertisement that tells you a certain drug will cure your particular ill. Do you run to your doctor and demand a prescription?



Although not all statistics are misleading and not everyone is out to get you, you do need to be vigilant. By sorting out the good information from the suspicious and bad information, you can steer clear of statistics that go wrong. The tools and strategies in this chapter are designed to help you to stop and say, "Wait a minute!" so you can analyze and critically think about the issues and make good decisions.

Chapter 4

Tools of the Trade

In This Chapter

- ▶ Seeing statistics as a process, not just as numbers
 - ▶ Getting familiar with some basic statistical jargon
-

In today's world, the buzzword is *data*, as in, "Do you have any data to support your claim?" "What data do you have on this?" "The data supported the original hypothesis that . . . , " "Statistical data show that . . . , " and "The data bear this out" But the field of statistics is not just about data.



Statistics is the entire process involved in gathering evidence to answer questions about the world, in cases where that evidence happens to be data.

In this chapter, you see firsthand how statistics works as a process and where the numbers play their part. You're also introduced to the most commonly used forms of statistical jargon, and you find out how these definitions and concepts all fit together as part of that process. So the next time you hear someone say, "This survey had a margin of error of plus or minus 3 percentage points," you'll have a basic idea of what that means.

Statistics: More than Just Numbers

Statisticians don't just "do statistics." Although the rest of the world views them as number crunchers, they think of themselves as the keepers of the scientific method. Of course, statisticians work with experts in other fields to satisfy their need for data, because man cannot live by statistics alone, but crunching someone's data is only a small part of a statistician's job. (In fact, if that's all we did all day, we'd quit our day jobs and moonlight as casino consultants.) In reality, statistics is involved in every aspect of the *scientific method* — formulating good questions, setting up studies, collecting good data, analyzing the data properly, and making appropriate conclusions. But aside from analyzing the data properly, what do any of these aspects have to do with statistics? In this chapter you find out.

All research starts with a question, such as:

- ✓ Is it possible to drink too much water?

- ✓ What's the cost of living in San Francisco?
- ✓ Who will win the next presidential election?
- ✓ Do herbs really help maintain good health?
- ✓ Will my favorite TV show get renewed for next year?

None of these questions asks anything directly about numbers. Yet each question requires the use of data and statistical processes to come up with the answer.

Suppose a researcher wants to determine who will win the next U.S. presidential election. To answer with confidence, the researcher has to follow several steps:

1. Determine the population to be studied.

In this case, the researcher intends to study registered voters who plan to vote in the next election.

2. Collect the data.

This step is a challenge, because you can't go out and ask every person in the United States whether they plan to vote, and if so, for whom they plan to vote. Beyond that, suppose someone says, "Yes, I plan to vote." Will that person *really* vote come Election Day? And will that same person tell you whom he actually plans to vote for? And what if that person changes his mind later on and votes for a different candidate?

3. Organize, summarize, and analyze the data.

After the researcher has gone out and collected the data she needs, getting it organized, summarized, and analyzed helps the researcher answer her question. This step is what most people recognize as the business of statistics.

4. Take all the data summaries, charts, graphs, and analyses and draw conclusions from them to try to answer the researcher's original question.

Of course, the researcher will not be able to have 100% confidence that her answer is correct, because not every person in the United States was asked. But she can get an answer that she is *nearly* 100% sure is correct. In fact, with a sample of about 2,500 people who are selected in a fair and *unbiased* way (that is, every possible sample of size 2,500 had an equal chance of being selected), the researcher can get accurate results within plus or minus 2.5% (if all the steps in the research process are done correctly).



In making conclusions, the researcher has to be aware that every study has limits and that — because the chance for error always exists — the results could be wrong. A numerical value can be reported that tells others how confident the researcher is about the results and how accurate these results are expected to be. (See Chapter 12 for more information on margin of error.)



After the research is done and the question has been answered, the results typically lead to even more questions and even more research. For example, if men appear to favor one candidate but women favor the opponent, the next questions may be: “Who goes to the polls more often on Election Day — men or women — and what factors determine whether they will vote?”

The field of statistics is really the business of using the scientific method to answer research questions about the world. Statistical methods are involved in every step of a good study, from designing the research to collecting the data, organizing and summarizing the information, doing an analysis, drawing conclusions, discussing limitations, and, finally, designing the next study in order to answer new questions that arise. Statistics is more than just numbers — it’s a process.

Grabbing Some Basic Statistical Jargon

Every trade has a basic set of tools, and statistics is no different. If you think about the statistical process as a series of stages that you go through to get from question to answer, you may guess that at each stage you’ll find a group of tools and a set of terms (or jargon) to go along with it. Now if the hair is beginning to stand up on the back of your neck, don’t worry. No one is asking you to become a statistics expert and plunge into the heavy-duty stuff, or to turn into a statistics nerd who uses this jargon all the time. Hey, you don’t even have to carry a calculator and pocket protector in your shirt pocket (because statisticians really don’t do that; it’s just an urban myth).

But as the world becomes more numbers-conscious, statistical terms are thrown around more in the media and in the workplace, so knowing what the language really means can give you a leg up. Also, if you’re reading this book because you want to find out more about how to calculate some statistics, understanding basic jargon is your first step. So, in this section, you get a taste of statistical jargon; I send you to the appropriate chapters elsewhere in the book to get details.

Data

Data are the actual pieces of information that you collect through your study. For example, I asked five of my friends how many pets they own, and the data they gave me are the following: 0, 2, 1, 4, 18. (The fifth friend counted each of her aquarium fish as a separate pet.) Not all data are numbers; I also recorded the gender of each of my friends, giving me the following data: male, male, female, male, female.

Most data fall into one of two groups: numerical or categorical. (I present the main ideas about these variables here; see Chapter 5 for more details.)

✓ **Numerical data:** These data have meaning as a measurement, such as a person's height, weight, IQ, or blood pressure; or they're a count, such as the number of stock shares a person owns, how many teeth a dog has, or how many pages you can read of your favorite book before you fall asleep. (Statisticians also call numerical data *quantitative data*.)

Numerical data can be further broken into two types: discrete and continuous.

- *Discrete data* represent items that can be counted; they take on possible values that can be listed out. The list of possible values may be fixed (also called *finite*); or it may go from 0, 1, 2, on to infinity (making it *countably infinite*). For example, the number of heads in 100 coin flips takes on values from 0 through 100 (finite case), but the number of flips needed to get 100 heads takes on values from 100 (the fastest scenario) on up to infinity. Its possible values are listed as 100, 101, 102, 103, . . . (representing the countably infinite case).
- *Continuous data* represent measurements; their possible values cannot be counted and can only be described using intervals on the real number line. For example, the exact amount of gas purchased at the pump for cars with 20-gallon tanks represents nearly-continuous data from 0.00 gallons to 20.00 gallons, represented by the interval [0, 20], inclusive. (Okay, you *can* count all these values, but why would you want to? In cases like these, statisticians bend the definition of continuous a wee bit.) The lifetime of a C battery can be anywhere from 0 to infinity, technically, with all possible values in between. Granted, you don't expect a battery to last more than a few hundred hours, but no one can put a cap on how long it can go (remember the Energizer Bunny?).

✓ **Categorical data:** Categorical data represent characteristics such as a person's gender, marital status, hometown, or the types of movies they like. Categorical data can take on numerical values (such as "1" indicating male and "2" indicating female), but those numbers don't have meaning. You couldn't add them together, for example. (Other names for categorical data are *qualitative data*, or *Yes/No data*.)



Ordinal data mixes numerical and categorical data. The data fall into categories, but the numbers placed on the categories have meaning. For example, rating a restaurant on a scale from 0 to 4 stars gives ordinal data. Ordinal data are often treated as categorical, where the groups are ordered when graphs and charts are made. I don't address them separately in this book.

Data set

A *data set* is the collection of all the data taken from your sample. For example, if you measured the weights of five packages, and those weights were 12, 15, 22, 68, and 3 pounds, those five numbers (12, 15, 22, 68, 3) constitute your data set. If you only record the general size of the package (for example, small, medium, or large), your data set may look like this: medium, medium, medium, large, small.

Variable

A *variable* is any characteristic or numerical value that varies from individual to individual. A variable can represent a count (for example, the number of pets you own); or a measurement (the time it takes you to wake up in the morning). Or the variable can be categorical, where each individual is placed into a group (or category) based on certain criteria (for example, political affiliation, race, or marital status). Actual pieces of information recorded on individuals regarding a variable are the data.

Population

For virtually any question you may want to investigate about the world, you have to center your attention on a particular group of individuals (for example, a group of people, cities, animals, rock specimens, exam scores, and so on). For example:

- ✓ What do Americans think about the president's foreign policy?
- ✓ What percentage of planted crops in Wisconsin did deer destroy last year?
- ✓ What's the prognosis for breast cancer patients taking a new experimental drug?
- ✓ What percentage of all cereal boxes get filled according to specification?

In each of these examples, a question is posed. And in each case, you can identify a specific group of individuals being studied: the American people, all planted crops in Wisconsin, all breast cancer patients, and all cereal boxes that are being filled, respectively. The group of individuals you want to study in order to answer your research question is called a *population*. Populations, however, can be hard to define. In a good study, researchers define the population very clearly, whereas in a bad study, the population is poorly defined.

The question of whether babies sleep better with music is a good example of how difficult defining the population can be. Exactly how would you define a baby? Under three months old? Under a year? And do you want to study babies only in the United States, or all babies worldwide? The results may be different for older and younger babies, for American versus European versus African babies, and so on.



Many times researchers want to study and make conclusions about a broad population, but in the end — to save time, money, or just because they don't know any better — they study only a narrowly defined population. That shortcut can lead to big trouble when conclusions are drawn. For example, suppose a college professor wants to study how TV ads persuade consumers to buy products. Her study is based on a group of her own students who participated to get five points extra credit. This test group may be convenient, but her results can't be generalized to any population beyond her own students, because no other population was represented in her study.

Sample, random, or otherwise

When you sample some soup, what do you do? You stir the pot, reach in with a spoon, take out a little bit of the soup, and taste it. Then you draw a conclusion about the whole pot of soup, without actually having tasted all of it. If your sample is taken in a fair way (for example, you didn't just grab all the good stuff) you will get a good idea how the soup tastes without having to eat it all. Taking a sample works the same way in statistics. Researchers want to find out something about a population, but they don't have time or money to study every single individual in the population. So they select a subset of individuals from the population, study those individuals, and use that information to draw conclusions about the whole population. This subset of the population is called a *sample*.

Although the idea of a selecting a sample seems straightforward, it's anything but. The way a sample is selected from the population can mean the difference between results that are correct and fair and results that are garbage. Example: Suppose you want a sample of teenagers' opinions on whether they're spending too much time on the Internet. If you send out a survey using text messaging, your results won't represent the opinions of *all teenagers*, which is your intended population. They will represent only those teenagers who have access to text messages. Does this sort of statistical mismatch happen often? You bet.



Some of the biggest culprits of statistical misrepresentation caused by bad sampling are surveys done on the Internet. You can find thousands of surveys on the Internet that are done by having people log on to a particular Web site and give their opinions. But even if 50,000 people in the U.S. complete a survey on the Internet, it doesn't represent the population of all Americans. It represents only those folks who have Internet access, who logged on to that particular Web site, and who were interested enough to participate in the survey (which typically means that they have strong opinions about the topic in question). The result of all these problems is *bias* — systematic favoritism of certain individuals or certain outcomes of the study.



How do you select a sample in a way that avoids bias? The key word is *random*. A *random sample* is a sample selected by equal opportunity; that is, every possible sample the same size as yours had an equal chance to be selected from the population. What *random* really means is that no group in the population is favored in or excluded from the selection process.

Non-random (in other words *bad*) *samples* are samples that were selected in such a way that some type of favoritism and/or automatic exclusion of a part of the population was involved. A classic example of a non-random sample comes from polls for which the media asks you to phone in your opinion on a certain issue (“call-in” polls). People who choose to participate in call-in polls do not represent the population at large because they had to be watching that program, and they had to feel strongly enough to call in. They technically don’t represent a sample at all, in the statistical sense of the word, because no one selected them beforehand — they selected themselves to participate, creating a *volunteer* or *self-selected* sample. The results will be skewed toward people with strong opinions.

To take an authentic random sample, you need a randomizing mechanism to select the individuals. For example, the Gallup Organization starts with a computerized list of all telephone exchanges in America, along with estimates of the number of residential households that have those exchanges. The computer uses a procedure called *random digit dialing* (RDD) to randomly create phone numbers from those exchanges, and then selects samples of telephone numbers from those. So what really happens is that the computer creates a list of *all possible* household phone numbers in America and then selects a subset of numbers from that list for Gallup to call.

Another example of random sampling involves the use of random number generators. In this process, the items in the sample are chosen using a computer-generated list of random numbers, where each sample of items has the same chance of being selected. Researchers may use this type of randomization to assign patients to a treatment group versus a control group in an experiment. This process is equivalent to drawing names out of a hat or drawing numbers in a lottery.



No matter how large a sample is, if it’s based on non-random methods, the results will not represent the population that the researcher wants to draw conclusions about. Don’t be taken in by large samples — first check to see how they were selected. Look for the term *random sample*. If you see that term, dig further into the fine print to see how the sample was actually selected and use the preceding definition to verify that the sample was, in fact, selected randomly. A small random sample is better than a large non-random one.

A *statistic* is a number that summarizes the data collected from a sample. People use many different statistics to summarize data. For example, data can be summarized as a percentage (60% of U.S. households sampled own more than two cars), an average (the average price of a home in this sample is . . .), a median (the median salary for the 1,000 computer scientists in this sample was . . .), or a percentile (your baby's weight is at the 90th percentile this month, based on data collected from over 10,000 babies).

The type of statistic calculated depends on the type of data. For example, percentages are used to summarize categorical data, and means are used to summarize numerical data. The price of a home is a numerical variable, so you can calculate its mean or standard deviation. However, the color of a home is a categorical variable; finding the standard deviation or median of color makes no sense. In this case, the important statistics are the percentages of homes of each color.



Not all statistics are correct or fair, of course. Just because someone gives you a statistic, nothing guarantees that the statistic is scientific or legitimate. You may have heard the saying, "Figures don't lie, but liars figure."

Parameter

Statistics are based on sample data, not on population data. If you collect data from the entire population, that process is called a *census*. If you then summarize the entire census information from one variable into a single number, that number is a *parameter*, not a statistic. Most of the time, researchers are trying to estimate the parameters using statistics. The U.S. Census Bureau wants to report the total number of people in the U.S., so it conducts a census. However, due to logistical problems in doing such an arduous task (such as being able to contact homeless folks), the census numbers can only be called *estimates* in the end, and they're adjusted upward to account for people the census missed.

Bias

Bias is a word you hear all the time, and you probably know that it means something bad. But what really constitutes bias? *Bias* is systematic favoritism that is present in the data collection process, resulting in lopsided, misleading results. Bias can occur in any of a number of ways:

- ✓ **In the way the sample is selected:** For example, if you want to estimate how much holiday shopping people in the United States plan to do this year, and you take your clipboard and head out to a shopping mall on the day after Thanksgiving to ask customers about their shopping plans, you have bias in your sampling process. Your sample tends to favor those die-hard shoppers at that

particular mall who were braving the massive crowds on that day known to retailers and shoppers as “Black Friday.”

- ✓ **In the way data are collected:** Poll questions are a major source of bias. Because researchers are often looking for a particular result, the questions they ask can often reflect and lead to that expected result. For example, the issue of a tax levy to help support local schools is something every voter faces at one time or another. A poll question asking, “Don’t you think it would be a great investment in our future to support the local schools?” has a bit of bias. On the other hand, so does “Aren’t you tired of paying money out of your pocket to educate other people’s children?” Question wording can have a huge impact on results.

Other issues that result in bias with polls are timing, length, level of question difficulty, and the manner in which the individuals in the sample were contacted (phone, mail, house-to-house, and so on). See Chapter 16 for more information on designing and evaluating polls and surveys.



When examining polling results that are important to you or that you’re particularly interested in, find out what questions were asked and exactly how the questions were worded before drawing your conclusions about the results.

Mean (Average)

The mean, also referred to by statisticians as the *average*, is the most common statistic used to measure the center, or middle, of a numerical data set. The *mean* is the sum of all the numbers divided by the total number of numbers. The mean of the entire population is called the *population mean*, and the mean of a sample is called the *sample mean*. (See Chapter 5 for more on the mean.)



The mean may not be a fair representation of the data, because the average is easily influenced by *outliers* (very small or large values in the data set that are not typical).

Median

The median is another way to measure the center of a numerical data set. A statistical median is much like the median of an interstate highway. On many highways, the median is the middle, and an equal number of lanes lay on either side of it. In a numerical data set, the *median* is the point at which there are an equal number of data points whose values lie above and below the median value. Thus, the median is truly the middle of the data set. See Chapter 5 for more on the median.



The next time you hear an average reported, look to see whether the median is also reported. If not, ask for it! The average and the median are two different representations of the middle of a data set and can often give two very different stories about the data, especially when the data set contains outliers (very large or small numbers that are not typical).

Standard deviation

Have you heard anyone report that a certain result was found to be “two standard deviations above the mean”? More and more, people want to report how significant their results are, and the number of standard deviations above or below average is one way to do it. But exactly what is a standard deviation?

The *standard deviation* is a measurement statisticians use for the amount of variability (or spread) among the numbers in a data set. As the term implies, a standard deviation is a standard (or typical) amount of deviation (or distance) from the average (or mean, as statisticians like to call it). So the standard deviation, in very rough terms, is the average distance from the mean.

The formula for standard deviation (denoted by s) is as follows, where n equals the number of values in the data set, each x represents a number in the data set, and \bar{x} is the average of all the data:

$$s = \sqrt{\sum \frac{(x - \bar{x})^2}{n-1}}$$

For detailed instructions on calculating the standard deviation, see Chapter 5.



The standard deviation is also used to describe where most of the data should fall, in a relative sense, compared to the average. For example, if your data have the form of a bell-shaped curve (also known as a *normal distribution*), about 95% of the data lie within two standard deviations of the mean. (This result is called the *empirical rule*, or the *68–95–99.7% rule*. See Chapter 5 for more on this.)



The standard deviation is an important statistic, but it is often absent when statistical results are reported. Without it, you’re getting only part of the story about the data. Statisticians like to tell the story about the man who had one foot in a bucket of ice water and the other foot in a bucket of boiling water. He said on average he felt just great! But think about the variability in the two temperatures for each of his feet. Closer to home, the average house price, for example, tells you nothing about the range of house prices you may encounter when house-hunting. The

average salary may not fully represent what's really going on in your company, if the salaries are extremely spread out.



Don't be satisfied with finding out only the average — be sure to ask for the standard deviation as well. Without a standard deviation, you have no way of knowing how spread out the values may be. (If you're talking starting salaries, for example, this could be very important!)

Percentile

You've probably heard references to percentiles before. If you've taken any kind of standardized test, you know that when your score was reported, it was presented to you with a measure of where you stood compared to the other people who took the test. This comparison measure was most likely reported to you in terms of a percentile. The *percentile* reported for a given score is the percentage of values in the data set that fall below that certain score. For example, if your score was reported to be at the 90th percentile, that means that 90% of the other people who took the test with you scored lower than you did (and 10% scored higher than you did). The median is right in the middle of a data set, so it represents the 50th percentile. For more specifics on percentiles, see Chapter 5.



Percentiles are used in a variety of ways for comparison purposes and to determine *relative standing* (that is, how an individual data value compares to the rest of the group). Babies' weights are often reported in terms of percentiles, for example. Percentiles are also used by companies to see where they stand compared to other companies in terms of sales, profits, customer satisfaction, and so on.

Standard score

The standard score is a slick way to put results in perspective without having to provide a lot of details — something that the media loves. The *standard score* represents the number of standard deviations above or below the mean (without caring what that standard deviation or mean actually are).

For example, suppose Bob took his statewide 10th-grade test recently and scored 400. What does that mean? Not much, because you can't put 400 into perspective. But knowing that Bob's standard score on the test is +2 tells you everything. It tells you that Bob's score is two standard deviations above the mean. (Bravo, Bob!) Now suppose Emily's standard score is -2. In this case, this is not good (for Emily), because it means her score is two standard deviations *below* the mean.

The process of taking a number and converting it to a standard score is called *standardizing*. For the details on calculating and interpreting standard scores when you have a normal (bell-shaped) distribution, see Chapter 9.

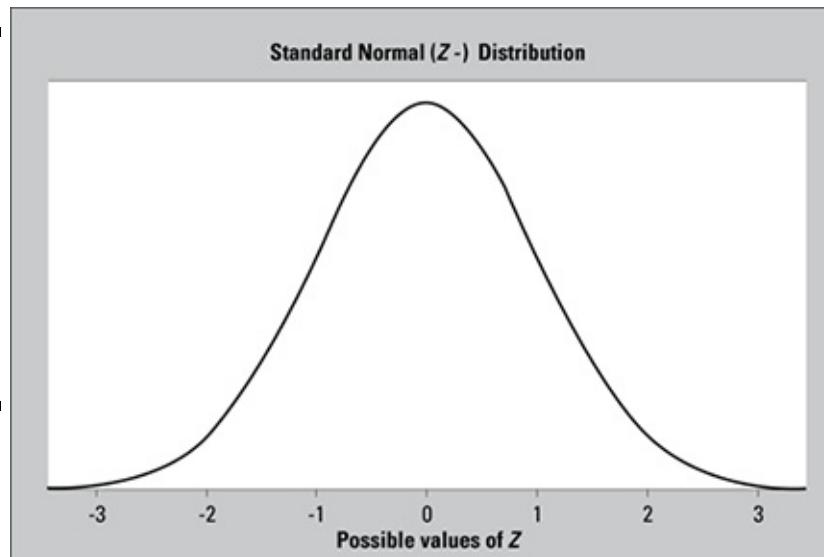
Distribution and normal distribution

The *distribution* of a data set (or a population) is a listing or function showing all the possible values (or intervals) of the data and how often they occur. When a distribution of categorical data is organized, you see the number or percentage of individuals in each group. When a distribution of numerical data is organized, they're often ordered from smallest to largest, broken into reasonably sized groups (if appropriate), and then put into graphs and charts to examine the shape, center, and amount of variability in the data.

The world of statistics includes dozens of different distributions for categorical and numerical data; the most common ones have their own names. One of the most well-known distributions is called the *normal distribution*, also known as the *bell-shaped curve*. The normal distribution is based on numerical data that is continuous; its possible values lie on the entire real number line. Its overall shape, when the data are organized in graph form, is a symmetric bell-shape. In other words, most (around 68%) of the data are centered around the mean (giving you the middle part of the bell), and as you move farther out on either side of the mean, you find fewer and fewer values (representing the downward sloping sides on either side of the bell).

The mean (and hence the median) is directly in the center of the normal distribution due to symmetry, and the standard deviation is measured by the distance from the mean to the *inflection point* (where the curvature of the bell changes from concave up to concave down). Figure 4-1 shows a graph of a normal distribution with mean 0 and standard deviation 1 (this distribution has a special name, the *standard normal distribution* or *Z-distribution*). The shape of the curve resembles the outline of a bell.

Figure 4-1: A standard normal (Z-) distribution has a bell-shaped curve with mean 0 and standard deviation 1.



Because every distinct population of data has a different mean and standard deviation, an infinite number of different normal distributions exist, each with its own mean and its own standard deviation to characterize it. See Chapter 9 for plenty more on the normal and standard normal distributions.

Central Limit Theorem



The normal distribution is also used to help measure the accuracy of many statistics, including the mean, using an important result in statistics called the *Central Limit Theorem*. This theorem gives you the ability to measure how much your sample mean will vary, without having to take any other sample means to compare it with (thankfully!). By taking this variability into account, you can now use your data to answer questions about the population, such as “What’s the mean household income for the whole U.S.”; or “This report said 75% of all gift cards go unused; is that really true?” (These two particular analyses made possible by the Central Limit Theorem are called *confidence intervals* and *hypothesis tests*, respectively, and are described in Chapters 13 and 14, respectively.)

The Central Limit Theorem (*CLT* for short) basically says that for non-normal data, your sample mean has an approximate normal distribution, no matter what the distribution of the original data looks like (as long as your sample size was large enough). And it doesn’t just apply to the sample mean; the CLT is also true for other sample statistics, such as the sample proportion (see Chapters 13 and 14). Because statisticians know so much about the normal distribution (see the preceding section), these analyses are much easier. See Chapter 11 for more on the Central Limit Theorem, known by statisticians as the “Crown jewel in the field of all statistics.” (Should you even bother to tell them to get a life?)

z-values



If a data set has a normal distribution, and you standardize all the data to obtain standard scores, those standard scores are called *z*-values. All *z*-values have what is known as a standard normal distribution (or *Z*-distribution). The *standard normal distribution* is a special normal distribution with a mean equal to 0 and a standard deviation equal to 1.

The standard normal distribution is useful for examining the data and determining statistics like percentiles, or the percentage of the data falling between two values. So if researchers determine that the data have a normal distribution, they usually first standardize the data (by converting each data point into a *z*-value) and then use the standard normal distribution to explore and discuss the data in more detail. See Chapter

9 for more details on *z*-values.

Experiments

An *experiment* is a study that imposes a treatment (or control) to the subjects (participants), controls their environment (for example, restricting their diets, giving them certain dosage levels of a drug or placebo, or asking them to stay awake for a prescribed period of time), and records the responses. The purpose of most experiments is to pinpoint a cause-and-effect relationship between two factors (such as alcohol consumption and impaired vision; or dosage level of a drug and intensity of side effects). Here are some typical questions that experiments try to answer:

- ✓ Does taking zinc help reduce the duration of a cold? Some studies show that it does.
- ✓ Does the shape and position of your pillow affect how well you sleep at night? The Emory Spine Center in Atlanta says yes.
- ✓ Does shoe heel height affect foot comfort? A study done at UCLA says up to one-inch heels are better than flat soles.

In this section, I discuss some additional definitions of words that you may hear when someone is talking about experiments. Chapter 17 is entirely dedicated to the subject. For now, just concentrate on basic experiment lingo.

Treatment group versus control group

Most experiments try to determine whether some type of experimental treatment (or important factor) has a significant effect on an outcome. For example, does zinc help to reduce the length of a cold? Subjects who are chosen to participate in the experiment are typically divided into two groups: a treatment group and a control group. (More than one treatment group is possible.)

- ✓ The *treatment group* consists of participants who receive the experimental treatment whose effect is being studied (in this case, zinc tablets).
- ✓ The *control group* consists of participants who do not receive the experimental treatment being studied. Instead, they get a placebo (a fake treatment; for example, a sugar pill); a standard, nonexperimental treatment (such as vitamin C, in the zinc study); or no treatment at all, depending on the situation.

In the end, the responses of those in the treatment group are compared with the responses from the control group to look for differences that are statistically significant (unlikely to have occurred just by chance).

Placebo

A *placebo* is a fake treatment, such as a sugar pill. Placebos are given to the control group to account for a psychological phenomenon called the *placebo effect*, in which patients receiving a fake treatment still report having a response, as if it were the real treatment. For example, after taking a sugar pill a patient experiencing the placebo effect might say, “Yes, I feel better already,” or “Wow, I *am* starting to feel a bit dizzy.” By measuring the placebo effect in the control group, you can tease out what portion of the reports from the treatment group were real and what portion were likely due to the placebo effect. (Experimenters assume that the placebo effect affects both the treatment and control groups.)

Blind and double-blind

A *blind experiment* is one in which the subjects who are participating in the study are not aware of whether they’re in the treatment group or the control group. In the zinc example, the vitamin C tablets and the zinc tablets would be made to look exactly alike and patients would not be told which type of pill they were taking. A blind experiment attempts to control for bias on the part of the participants.

A *double-blind experiment* controls for potential bias on the part of both the patients *and* the researchers. Neither the patients nor the researchers collecting the data know which subjects received the treatment and which didn’t. So who does know what’s going on as far as who gets what treatment? Typically a third party (someone not otherwise involved in the experiment) puts together the pieces independently. A double-blind study is best, because even though researchers may claim to be unbiased, they often have a special interest in the results — otherwise they wouldn’t be doing the study!

Surveys (Polls)

A *survey* (more commonly known as a *poll*) is a questionnaire; it’s most often used to gather people’s opinions along with some relevant demographic information. Because so many policymakers, marketers, and others want to “get at the pulse of the American public” and find out what the average American is thinking and feeling, many people now feel that they cannot escape the barrage of requests to take part in surveys and polls. In fact, you’ve probably received many requests to participate in surveys, and you may even have become numb to them, simply throwing away surveys received in the mail or saying “no” when asked to participate in a telephone survey.

If done properly, a good survey can really be informative. People use surveys to find out what TV programs Americans (and others) like, how consumers feel about Internet shopping, and whether the United States should allow someone under 35 to become president. Surveys are used by companies to assess the level of satisfaction their customers feel, to find out what products their customers want, and to determine who is

buying their products. TV stations use surveys to get instant reactions to news stories and events, and movie producers use them to determine how to end their movies.

However, if I had to choose one word to assess the general state of surveys in the media today, I'd say it's *quantity* rather than *quality*. In other words, you'll find no shortage of bad surveys. But in this book you find no shortage of good tips and information for analyzing, critiquing, and understanding survey results, and for designing your own surveys to do the job right. (To take off with surveys, head to Chapter 16.)

Margin of error

You've probably heard or seen results like this: "This survey had a margin of error of plus or minus 3 percentage points." What does this mean? Most surveys (except a census) are based on information collected from a sample of individuals, not the entire population. A certain amount of error is bound to occur — not in the sense of calculation error (although there may be some of that, too) but in the sense of *sampling error*, which is the error that occurs simply because the researchers aren't asking everyone. The *margin of error* is supposed to measure the maximum amount by which the sample results are expected to differ from those of the actual population. Because the results of most survey questions can be reported in terms of percentages, the margin of error most often appears as a percentage, as well.

How do you interpret a margin of error? Suppose you know that 51% of people sampled say that they plan to vote for Ms. Calculation in the upcoming election. Now, projecting these results to the whole voting population, you would have to add and subtract the margin of error and give a range of possible results in order to have sufficient confidence that you're bridging the gap between your sample and the population. Supposing a margin of error of plus or minus 3 percentage points, you would be pretty confident that between 48% ($51\% - 3\%$) and 54% ($51\% + 3\%$) of the population will vote for Ms. Calculation in the election, based on the sample results. In this case, Ms. Calculation may get slightly more or slightly less than the majority of votes and could either win or lose the election. This has become a familiar situation in recent years when the media want to report results on Election Night, but based on early exit polling results, the election is "too close to call." For more on the margin of error, see Chapter 12.



The margin of error measures accuracy; it does not measure the amount of bias that may be present (find a discussion of bias earlier in this chapter). Results that look numerically scientific and precise don't mean anything if they were collected in a biased way.

Confidence interval

One of the biggest uses of statistics is to estimate a population parameter using a sample statistic. In other words, use a number that summarizes a sample to help you guesstimate the corresponding number that summarizes the whole population (the definitions of parameter and statistic appear earlier in this chapter). You’re looking for a population parameter in each of the following questions:

- ✓ What’s the average household income in America? (Population = all households in America; parameter = average household income.)
- ✓ What percentage of all Americans watched the Academy Awards this year? (Population = all Americans; parameter = percentage who watched the Academy Awards this year.)
- ✓ What’s the average life expectancy of a baby born today? (Population = all babies born today; parameter = average life expectancy.)
- ✓ How effective is this new drug on adults with Alzheimer’s? (Population = all people who have Alzheimer’s; parameter = percentage of these people who see improvement when taking this drug.)

It’s not possible to find these parameters exactly; they each require an estimate based on a sample. You start by taking a random sample from a population (say a sample of 1,000 households in America) and then finding the corresponding statistic from that sample (the sample’s mean household income). Because you know that sample results vary from sample to sample, you need to add a “plus or minus something” to your sample results if you want to draw conclusions about the whole population (all households in America). This “plus or minus” that you add to your sample statistic in order to estimate a parameter is the margin of error.

When you take a sample statistic (such as the sample mean or sample percentage) and add/subtract a margin of error, you come up with what statisticians call a *confidence interval*. A confidence interval represents a range of likely values for the population parameter, based on your sample statistic. For example, suppose the average time it takes you to drive to work each day is 35 minutes, with a margin of error of plus or minus 5 minutes. You estimate that the average time to work would be anywhere from 30 to 40 minutes. This estimate is a confidence interval.



Some confidence intervals are wider than others (and wide isn’t good, because it equals less accuracy). Several factors influence the width of a confidence interval, such as sample size, the amount of variability in the population being studied, and how confident you want to be in your results. (Most researchers are happy with a 95% level of confidence in their results.) For more on factors that influence confidence intervals, as well as instructions for calculating and interpreting confidence intervals, see Chapter 13.

Hypothesis testing

Hypothesis test is a term you probably haven't run across in your everyday dealings with numbers and statistics. But I guarantee that hypothesis tests have been a big part of your life and your workplace, simply because of the major role they play in industry, medicine, agriculture, government, and a host of other areas. Any time you hear someone talking about their study showing a "statistically significant result," you're encountering a hypothesis test. (A statistically significant result is one that is unlikely to have occurred by chance. See Chapter 14 for the full scoop.)

Basically, a *hypothesis test* is a statistical procedure in which data are collected from a sample and measured against a claim about a population parameter. For example, if a pizza delivery chain claims to deliver all pizzas within 30 minutes of placing the order, on average, you could test whether this claim is true by collecting a random sample of delivery times over a certain period and looking at the average delivery time for that sample. To make your decision, you must also take into account the amount by which your sample results can change from sample to sample (which is related to the margin of error).



Because your decision is based on a sample and not the entire population, a hypothesis test can sometimes lead you to the wrong conclusion. However, statistics are all you have, and if done properly, they can give you a good chance of being correct. For more on the basics of hypothesis testing, see Chapter 14.

A variety of hypothesis tests are done in scientific research, including *t*-tests (comparing two population means), paired *t*-tests (looking at before/after data), and tests of claims made about proportions or means for one or more populations. For specifics on these hypothesis tests, see Chapter 15.

p-values

Hypothesis tests are used to test the validity of a claim that is made about a population. This claim that's on trial, in essence, is called the *null hypothesis*. The *alternative hypothesis* is the one you would believe if the null hypothesis is concluded to be untrue. The evidence in the trial is your data and the statistics that go along with it. All hypothesis tests ultimately use a *p*-value to weigh the strength of the evidence (what the data are telling you about the population). The *p*-value is a number between 0 and 1 and interpreted in the following way:

- ✓ A small *p*-value (≤ 0.05) indicates strong evidence against the null hypothesis, so you reject it.
- ✓ A large *p*-value (> 0.05) indicates weak evidence against the null hypothesis, so

you fail to reject it.

- ✓ *p*-values very close to the cutoff (0.05) are considered to be marginal (could go either way). Always report the *p*-value so your readers can draw their own conclusions.

For example, suppose a pizza place claims their delivery times are 30 minutes or less on average but you think it's more than that. You conduct a hypothesis test because you believe the null hypothesis, H_0 , that the mean delivery time is 30 minutes max, is incorrect. Your alternative hypothesis (H_a) is that the mean time is greater than 30 minutes. You randomly sample some delivery times and run the data through the hypothesis test, and your *p*-value turns out to be 0.001, which is much less than 0.05. You conclude that the pizza place is wrong; their delivery times are in fact more than 30 minutes on average, and you want to know what they're gonna do about it! (Of course you could be wrong by having sampled an unusually high number of late pizzas just by chance; but whose side am I on?) For more on *p*-values, head to Chapter 14.

Statistical significance

Whenever data are collected to perform a hypothesis test, the researcher is typically looking for something out of the ordinary. (Unfortunately, research that simply confirms something that was already well known doesn't make headlines.) Statisticians measure the amount by which a result is out of the ordinary using hypothesis tests (see Chapter 14). They define a *statistically significant* result as a result with a very small probability of happening just by chance, and provide a number called a *p*-value to reflect that probability (see the previous section on *p*-values).

For example, if a drug is found to be more effective at treating breast cancer than the current treatment is, researchers say that the new drug shows a statistically significant improvement in the survival rate of patients with breast cancer. That means that based on their data, the difference in the overall results from patients on the new drug compared to those using the old treatment is so big that it would be hard to say it was just a coincidence. However, proceed with caution: You can't say that these results necessarily apply to each individual or to each individual in the same way. For full details on statistical significance, see Chapter 14.



When you hear that a study's results are statistically significant, don't automatically assume that the study's results are important. *Statistically significant* means the results were unusual, but unusual doesn't always mean important. For example, would you be excited to learn that cats move their tails more often when lying in the sun than when lying in the shade, and that those results are statistically significant? This result may not even be important to the cat, much less anyone else!

Sometimes statisticians make the wrong conclusion about the null hypothesis because a sample doesn't represent the population (just by chance). For example, a positive effect that's experienced by a sample of people who took the new treatment may have just been a fluke; or in the example in the preceding section, the pizza company really was delivering those pizzas on time and you just got an unlucky sample of slow ones. However, the beauty of research is that as soon as someone gives a press release saying that she found something significant, the rush is on to try to replicate the results, and if the results can't be replicated, this probably means that the original results were wrong for some reason (including being wrong just by chance). Unfortunately, a press release announcing a "major breakthrough" tends to get a lot of play in the media, but follow-up studies refuting those results often don't show up on the front page.



One statistically significant result shouldn't lead to quick decisions on anyone's part. In science, what most often counts is not a single remarkable study, but a body of evidence that is built up over time, along with a variety of well-designed follow-up studies. Take any major breakthroughs you hear about with a grain of salt and wait until the follow-up work has been done before using the information from a single study to make important decisions in your life. The results may not be replicable, and even if they are, you can't know if they necessarily apply to each individual.

Correlation versus causation



Of all of the misunderstood statistical issues, the one that's perhaps the most problematic is the misuse of the concepts of correlation and causation.

Correlation, as a statistical term, is the extent to which two numerical variables have a linear relationship (that is, a relationship that increases or decreases at a constant rate). Following are three examples of correlated variables:

- ✓ The number of times a cricket chirps per second is strongly related to temperature; when it's cold outside, they chirp less frequently, and as the temperature warms up, they chirp at a steadily increasing rate. In statistical terms, you say number of cricket chirps and temperature have a strong positive correlation.
- ✓ The number of crimes (per capita) has often been found to be related to the number of police officers in a given area. When more police officers patrol the area, crime tends to be lower, and when fewer police officers are present in the same area, crime tends to be higher. In statistical terms we say the number of police officers and the number of crimes have a strong negative correlation.
- ✓ The consumption of ice cream (pints per person) and the number of murders in

New York are positively correlated. That is, as the amount of ice cream sold per person increases, the number of murders increases. Strange but true!

But correlation as a statistic isn't able to explain *why* or *how* the relationship between two variables, x and y , exists; only that it does exist.

Causation goes a step further than correlation, stating that a change in the value of the x variable *will cause* a change in the value of the y variable. Too many times in research, in the media, or in the public consumption of statistical results, that leap is made when it shouldn't be. For instance, you can't claim that consumption of ice cream *causes* an increase in murder rates just because they are correlated. In fact, the study showed that temperature was positively correlated with both ice cream sales and murders. (For more on correlation and causation, see Chapter 18.) When can you make the causation leap? The most compelling case is when a well-designed experiment is conducted that rules out other factors that could be related to the outcomes (see Chapter 17 for information on experiments showing cause-and-effect).



You may find yourself wanting to jump to a cause-and-effect relationship when a correlation is found; researchers, the media, and the general public do it all the time. However, before making any conclusions, look at how the data were collected and/or wait to see if other researchers are able to replicate the results (the first thing they try to do after someone else's "groundbreaking result" hits the airwaves).

Part II

Number-Crunching Basics

The 5th Wave

By Rich Tennant



*"I ran an evaluation of our last pie chart.
Apparently it's boysenberry."*

In this part . . .

Number crunching: It's a dirty job, but somebody has to do it. Why not let it be you? Even if you aren't a numbers person and calculations aren't your thing, the step-by-step approach in this part may be just what you need to boost your confidence in doing and really understanding statistics.

In this part, you get down to the basics of number crunching, from making and interpreting charts and graphs to cranking out and understanding means, medians, standard deviations, and more. You also develop important skills for critiquing someone else's statistical information and getting at the real truth behind the data.

Chapter 5

Means, Medians, and More

In This Chapter

- ▶ Summarizing data effectively
 - ▶ Interpreting commonly used statistics
 - ▶ Realizing what statistics do and don't say
-

Every data set has a story, and if statistics are used properly, they do a good job of uncovering and reporting that story. Statistics that are improperly used can tell a different story, or only part of it, so knowing how to make good decisions about the information you're given is very important.

A *descriptive statistic* (or *statistic* for short) is a number that summarizes or describes some characteristic about a set of data. In this chapter, you see some of the most common descriptive statistics and how they are used, and you find out how to calculate them, interpret them, and put them together to get a good picture of a data set. You also find out what these statistics say and what they don't say about the data.

Summing Up Data with Descriptive Statistics

Descriptive statistics take a data set and boil it down to a set of basic information. Summarized data are often used to provide people with information that is easy to understand and that helps answer their questions. Picture your boss coming to you and asking, “What’s our client base like these days, and who’s buying our products?” How would you like to answer that question — with a long, detailed, and complicated stream of numbers that are sure to glaze her eyes over? Probably not. You want clean, clear, and concise statistics that sum up the client base for her, so that she can see how brilliant you are and then send you off to collect even more data to see how she can include more people in the client base. (That’s what you get for being efficient.)

Summarizing data has other purposes, as well. After all the data have been collected from a survey or some other kind of study, the next step is for the researcher to try to make sense out of the data. Typically, the first step researchers take is to run some basic statistics on the data to get a rough idea about what’s happening in it. Later in the process, researchers can do more analyses to formulate or test claims made about the population the data came from, estimate certain characteristics about the population (like the mean), look for links between variables they measured, and so on.

Another big part of research is reporting the results, not only to your peers, but also to the media and the general public. Although a researcher's peers may be anxiously waiting to hear about all the complex analyses that were done on a data set, the general public is neither ready for nor interested in that. What does the public want? Basic information. Statistics that make a point clearly and concisely are usually used to relay information to the media and to the public.



If you really need to learn more from data, a quick statistical overview isn't enough. In the statistical world, less is not more, and sometimes the real story behind the data can get lost in the shuffle. To be an informed consumer of statistics, you need to think about which statistics are being reported, what these statistics really mean, and what information is missing. This chapter focuses on these issues.

Crunching Categorical Data: Tables and Percents

Categorical data (also known as *qualitative data*) capture qualities or characteristics about the individual, such as a person's eye color, gender, political party, or opinion on some issue (using categories such as Agree, Disagree, or No opinion). Categorical data tend to fall into groups or categories pretty naturally. "Political party," for example, typically has four groups in the United States: Democrat, Republican, Independent, and Other. Categorical data often come from survey data, but they can also be collected in experiments. For example, in an experimental test of a new medical treatment, researchers may use three categories to assess the outcome of the experiment: Did the patient get better, worse, or stay the same while undergoing the treatment?

Categorical data are often summarized by reporting the percentage of individuals falling into each category. For example, pollsters may report political affiliation statistics by giving the percentage of Republicans, Democrats, Independents, and Others. To calculate the percentage of individuals in a certain category, find the number of individuals in that category, divide by the total number of people in the study, and then multiply by 100%. For example, if a survey of 2,000 teenagers included 1,200 females and 800 males, the resulting percentages would be $(1,200 \div 2,000) * 100\% = 60\%$ female and $(800 \div 2,000) * 100\% = 40\%$ male.

You can break down categorical data further by creating something called two-way tables. *Two-way tables* (also called *crosstabs*) are tables with rows and columns. They summarize the information from two categorical variables at once, such as gender and political party, so you can see (or easily calculate) the percentage of individuals in each combination of categories and use them to make comparisons between groups.

For example, if you had data about the gender and political party of your respondents, you would be able to look at the percentage of Republican females, Republican males, Democratic females, Democratic males, and so on. In this example, the total number of possible combinations in your table would be $2 * 4 = 8$, or the total number of gender categories times the total number of party affiliation categories. (See Chapter 19 for the full scoop, and then some, on two-way tables.)

The U.S. government calculates and summarizes loads of categorical data using crosstabs. Typical age and gender data, reported by the U.S. Census Bureau for a survey conducted in 2009, are shown in Table 5-1. (Normally age would be considered a numerical variable, but the way the U.S. government reports it, age is broken down into categories, making it a categorical variable.)

Table 5-1 **U.S. Population, Broken Down by Age and Gender (2009)**

Age Group	Both Sexes	%	Males	%	Females	%
Under 5	21,299,656	6.94	10,887,008	7.19	10,412,648	6.69
5–9	20,609,634	6.71	10,535,900	6.96	10,073,734	6.48
10–14	19,973,564	6.51	10,222,522	6.75	9,751,042	6.27
15–19	21,537,837	7.02	11,051,289	7.30	10,486,548	6.74
20–24	21,539,559	7.02	11,093,552	7.32	10,446,007	6.72
25–29	21,677,719	7.06	11,115,560	7.34	10,562,159	6.79
30–34	19,888,603	6.48	10,107,974	6.67	9,780,629	6.29
35–39	20,538,351	6.69	10,353,016	6.84	10,185,335	6.55
40–44	20,991,605	6.84	10,504,139	6.94	10,487,466	6.74

Table 5-1

Age Group	Both Sexes	%	Males	%	Females	%
45–49	22,831,092	7.44	11,295,524	7.46	11,535,568	7.42
50–54	21,761,391	7.09	10,677,847	7.05	11,083,544	7.13
55–59	18,975,026	6.18	9,204,666	6.08	9,770,360	6.28
60–64	15,811,923	5.15	7,576,933	5.00	8,234,990	5.29
65–69	11,784,320	3.84	5,511,164	3.64	6,273,156	4.03
70–74	9,007,747	2.93	4,082,226	2.70	4,925,521	3.17
75–79	7,325,528	2.39	3,149,236	2.08	4,176,292	2.68
80–84	5,822,334	1.90	2,298,260	1.52	3,524,074	2.27
85–89	3,662,397	1.19	1,266,899	0.84	2,395,498	1.54
90–94	1,502,263	0.49	424,882	0.28	1,077,381	0.69
95–99	401,977	0.13	82,135	0.05	319,842	0.21
100+	64,024	0.02	8,758	0.01	55,266	0.04
Total	307,006,550	100.00	151,449,490	100.00	155,557,060	100.00

You can examine many different facets of the U.S. population by looking at and working with different numbers from Table 5-1. For example, looking at gender, you notice that women slightly outnumber men — the population in 2009 was 50.67% female (divide total number of females by total population size and multiply by 100%) and 49.33% male (divide total number of males by total population size and multiply by 100%). You can also look at age: The percentage of the entire population that is under 5 years old was 6.94% (divide the total number under age 5 by the total population size and multiply by 100%). The largest group belongs to the 45–49 year olds, who made up 7.44% of the population.

Next, you can explore a possible relationship between gender and age by comparing various parts of the table. You can compare, for example, the percentage of females to males in the 80-and-over age group. Because these data are reported in 5-year increments, you have to do a little math in order to get your answer, though. The percentage of the population that's female and aged 80 and above (looking at column 7 of Table 5-1) is $2.27\% + 1.54\% + 0.69\% + 0.21\% + 0.04\% = 4.75\%$. The percentage of males aged 80 and over (looking at column 5 of Table 5-1) is $1.52\% + 0.84\% + 0.28\% + 0.05\% + 0.01\% = 2.70\%$. This shows that the 80-and-over age group for the females is about 76% larger than the males (because $[4.75 - 2.70] \div 2.70 = 0.76$).

These data confirm the widely accepted notion that women tend to live longer than men. However, the gap between men and women is narrowing over time. According to the U.S. Census Bureau, back in 2001 the percentage of women who were 80 years old and over was 4.36, compared to 2.31 for the men. The females in this age group outnumbered the males by a whopping 89% back in 2001 (note that $[4.36 - 2.31] \div 2.31 = 0.89$).



After you have the crosstabs that show the breakdown of two categorical variables, you can conduct hypothesis tests to determine whether a significant relationship or link between the two variables exists, taking into account the fact that data vary from sample to sample. Chapter 14 gives you all the details on hypothesis tests.

Measuring the Center with Mean and Median

With *numerical data*, measurable characteristics such as height, weight, IQ, age, or income are represented by numbers that make sense within the context of the problem (for example in units of feet, dollars, or people). Because the data have numerical meaning, you can summarize them in more ways than is possible with categorical data. The most common way to summarize a numerical data set is to describe where the center is. One way of thinking about what the center of a data set means is to ask, “What’s a typical value?” Or, “Where is the middle of the data?” The center of a data set can actually be measured in different ways, and the method chosen can greatly influence the

conclusions people make about the data. This section hits on measures of center.

Averaging out to the mean

NBA players make a lot of money, right? You often hear about players like Kobe Bryant or LeBron James who make tens of millions of dollars a year. But is that what the typical NBA player makes? Not really (although I don't exactly feel sorry for the others, given that they still make more money than most of us will ever make). Tens of millions of dollars is the kind of money you can command when you are a superstar among superstars, which is what these elite players are.

So how much money does the typical NBA player make? One way to answer this is to look at the average (the most commonly used statistic of all time).

The *average*, also called the *mean* of a data set, is denoted \bar{x} . The formula for finding the mean is:

$$\bar{x} = \frac{\sum x_i}{n}$$

where each value in the data set is denoted by an x with a subscript i that goes from 1 (the first number) to n (the last number).

Here's how you calculate the mean of a data set:

- 1. Add up all the numbers in the data set.**
- 2. Divide by the number of numbers in the data set, n .**



The mean I discuss here applies to a sample of data and is technically called the *sample mean*. The mean of an entire population of data is denoted with the Greek letter μ and is called the *population mean*. It's found by summing up all the values in the population and dividing by the population size, denoted N (to distinguish it from a sample size, n). Typically the population mean is unknown, and you use a sample mean to estimate it (plus or minus a margin of error; see all the details in Chapter 13).

For example, player salary data for the 13 players on the 2010 NBA Champion Los Angeles Lakers is shown in Table 5-2.

Table 5-2
**Salaries for L.A. Lakers
NBA Players (2009–2010)**

Player	Salary (\$)
Kobe Bryant	23,034,375
Pau Gasol	16,452,000
Andrew Bynum	12,526,998
Lamar Odom	7,500,000
Ron Artest	5,854,000
Adam Morrison	5,257,229
Derek Fisher	5,048,000
Sasha Vujacic	5,000,000
Luke Walton	4,840,000

Player	Salary (\$)
Shannon Brown	2,000,000
Jordan Farmar	1,947,240
Didier Ilunga-Mbenga	959,111
Josh Powell	959,111
Total	91,378,064

The mean of all the salaries on this team is $\$91,378,064 \div 13 = \$7,029,082$. That's a pretty nice average salary, isn't it? But notice that Kobe Bryant really stands out at the top of this list, and he should — his salary was the second highest in the entire league that season (just behind Tracy McGrady). If you remove Kobe from the equation (literally), the average salary of all the Lakers players besides Kobe becomes $\$68,343,689 \div 12 = \$5,695,307$ — a difference of around 1.3 million.

This new mean is still a hefty amount, but it's significantly lower than the mean salary of all the players including Kobe. (Fans would tell you that this reflects his importance to the team, and others would say no one is worth that much money; this issue is but the tip of the iceberg of the never-ending debates that sports fans — me included — love to have about statistics.)

Bottom line: The mean doesn't always tell the whole story. In some cases it may be a bit misleading, and this is one of those cases. That's because every year a few top-notch players (like Kobe) make much more money than anybody else, and their salaries pull up the overall average salary.



Numbers in a data set that are extremely high or extremely low compared to the rest of the data are called *outliers*. Because of the way the average is calculated, high outliers tend to drive the average upward (as Kobe's salary did in the preceding example). Low outliers tend to drive the average downward.

Splitting your data down the median

Remember in school when you took an exam, and you and most of the rest of the class did badly, but a couple of nerds got 100? Remember how the teacher didn't curve the scores to reflect the poor performance of most of the class? Your teacher was probably using the average, and the average in that case didn't really represent what statisticians might consider the best measure of center for the students' scores.

What can you report, other than the average, to show what the salary of a "typical" NBA player would be or what the test score of a "typical" student in your class was? Another statistic used to measure the center of a data set is called the median. The median is still an unsung hero of statistics in the sense that it isn't used nearly as often as it should be, although people are beginning to report it more nowadays.

The *median* of a data set is the value that lies exactly in the middle when the data have been ordered. It's denoted in different ways; some people use M and some use \bar{x} . Here are the steps for finding the median of a data set:

- 1. Order the numbers from smallest to largest.**
- 2. If the data set contains an odd number of numbers, choose the one that is exactly in the middle. You've found the median.**
- 3. If the data set contains an even number of numbers, take the two numbers that appear in the middle and average them to find the median.**

The salaries for the Los Angeles Lakers during the 2009–2010 season (refer to Table 5-2) are ordered from smallest (at the bottom) to largest (at the top). Because the list contains the names and salaries of 13 players, the middle salary is the seventh one from the bottom: Derek Fisher, who earned \$5.048 million that season from the Lakers. Derek is at the median.



This median salary (\$5.048 million) is well below the average of \$7.029 million for the 2009–2010 Lakers team. Notice that only 4 players of the 13 earned more than the average Lakers salary of \$7.029 million. Because the average includes outliers (like the salary of Kobe Bryant), the median salary is more representative of center for the team salaries. The median isn't affected by the salaries of those players who are way out there on the high end the way the average is.

Note: By the way, the lowest Lakers' salary for the 2009–2010 season was \$959,111 — a lot of money by most people's standards, but peanuts compared to what you imagine when you think of an NBA player's salary!



The U.S. government most often uses the median to represent the center with

respect to income data again because the median is not affected by outliers. For example, the U.S. Census Bureau reported that in 2008, the median household income was \$50,233 while the mean was found to be \$68,424. That's quite a difference!

Comparing means and medians: Histograms

Sometimes the mean versus median debate can get quite interesting. Suppose you're part of an NBA team trying to negotiate salaries. If you represent the owners, you want to show how much everyone is making and how much money you're spending, so you want to take into account those superstar players and report the average. But if you're on the side of the players, you would want to report the median, because that's more representative of what the players in the middle are making. Fifty percent of the players make a salary above the median, and 50 percent make a salary below the median. To sort it all out, it's best to find and compare both the mean and the median. A graph showing the shape of the data is a great place to start.



One of the graphs you can make to illustrate the shape of numerical data (how many values are close to/far from the mean, where the center is, how many outliers there might be) is a histogram. A *histogram* is a graph that organizes and displays numerical data in picture form, showing groups of data and the number or percentage of the data that fall into each group. It gives you a nice snapshot of the data set. (See Chapter 7 for more information on histograms and other types of data displays.)

Data sets can have many different possible shapes; here is a sampling of three shapes that are commonly discussed in introductory statistics courses:

- ✓ If most of the data are on the left side of the histogram but a few larger values are on the right, the data are said to be *skewed to the right*.

Histogram A in Figure 5-1 shows an example of data that are skewed to the right. The few larger values bring the mean upwards but don't really affect the median. So when data are skewed right, *the mean is larger than the median*. An example of such data is NBA salaries.

- ✓ If most of the data are on the right, with a few smaller values showing up on the left side of the histogram, the data are *skewed to the left*.

Histogram B in Figure 5-1 shows an example of data that are skewed to the left. The few smaller values bring the mean down, and again the median is minimally affected (if at all). An example of skewed-left data is the amount of time students use to take an exam; some students leave early, more of them stay later, and many stay until the bitter end (some would stay forever if they could!). When data are

skewed left, the mean is smaller than the median.

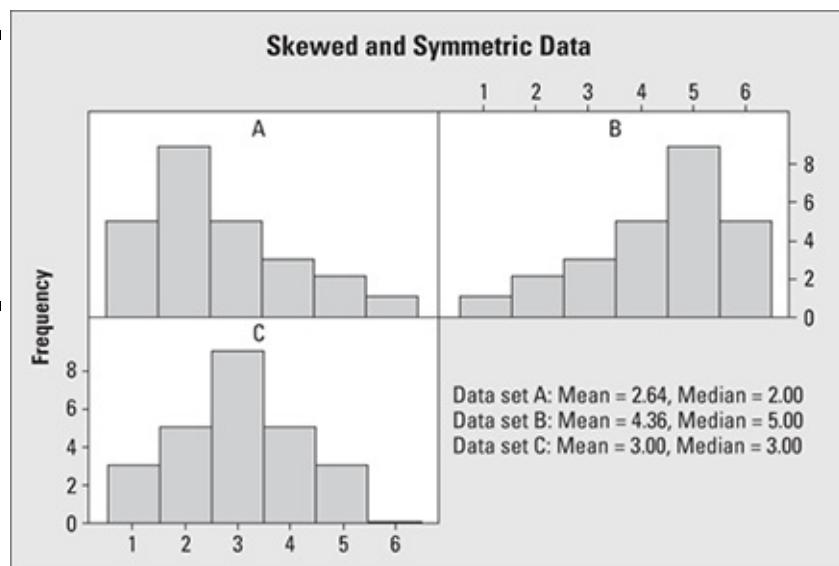
- ✓ If the data are *symmetric*, they have about the same shape on either side of the middle. In other words, if you fold the histogram in half, it looks about the same on both sides.

Histogram C in Figure 5-1 shows an example of symmetric data in a histogram. With symmetric data, the mean and median are close together.



By looking at Histogram A in Figure 5-1 (whose shape is skewed right), you can see that the “tail” of the graph (where the bars are getting shorter) is to the right, while the “tail” is to the left in Histogram B (whose shape is skewed left). By looking at the direction of the tail of a skewed distribution, you determine the direction of the skewness. Always add the direction when describing a skewed distribution.

Figure 5-1: A)
Data skewed
right; B) data
skewed left;
and C)
symmetric
data.



Histogram C is symmetric (it has about the same shape on each side). However, not all symmetric data has a bell shape like Histogram C does. As long as the shape is approximately the same on both sides, then you say that the shape is symmetric.



The average (or mean) of a data set is affected by outliers, but the median is not. In statistical lingo, if a statistic is not affected by a certain characteristic of the data (such as outliers, or skewness), then you say that statistic is *resistant* to that characteristic. In this case the median is resistant to outliers; the mean is not. If someone reports the average value, also ask for the median so that you can compare the two statistics and get a better feel for what's actually going on in the data and what's truly typical.

Accounting for Variation

Variation always exists in a data set, regardless of which characteristics you're measuring, because not every individual is going to have the same exact value for every variable. Variation is what makes the field of statistics what it is. For example, the price of homes varies from house to house, from year to year, and from state to state. The amount of time it takes you to get to work varies from day to day. The trick to dealing with variation is to be able to measure that variation in a way that best captures it.

Reporting the standard deviation

By far the most common measure of variation for numerical data is the standard deviation. The *standard deviation* measures how concentrated the data are around the mean; the more concentrated, the smaller the standard deviation. It's not reported nearly as often as it should be, but when it is, you often see it in parentheses: ($s = 2.68$).

Calculating standard deviation

The formula for the sample standard deviation of a data set (s) is

$$s = \sqrt{\frac{\sum(x - \bar{x})^2}{n-1}}$$

To calculate s , do the following steps:

1. Find the average of the data set, \bar{x} .
2. Take each number in the data set (x) and subtract the mean from it to get $(x - \bar{x})$.
3. Square each of the differences, $(x - \bar{x})^2$.
4. Add up all of the results from Step 3 to get the sum of squares: $\sum(x - \bar{x})^2$.
5. Divide the sum of squares (found in Step 4) by the number of numbers in the data set minus one; that is, $(n - 1)$. Now you have:

$$\frac{\sum(x - \bar{x})^2}{n-1}$$

6. Take the square root to get

$$s = \sqrt{\frac{\sum(x - \bar{x})^2}{n-1}}$$

which is the sample standard deviation, s . Whew!



At the end of Step 5 you have found a statistic called the *sample variance*, denoted by s^2 . The variance is another way to measure variation in a data set; its downside is that it's in square units. If your data are in dollars, for example, the variance would be in square dollars — which makes no sense. That's why we proceed to Step 6. Standard deviation has the same units as the original data.

Look at the following small example: Suppose you have four quiz scores: 1, 3, 5, and 7. The mean is $16 \div 4 = 4$ points. Subtracting the mean from each number, you get $(1 - 4) = -3$, $(3 - 4) = -1$, $(5 - 4) = +1$, and $(7 - 4) = +3$. Squaring each of these results, you get 9, 1, 1, and 9. Adding these up, the sum is 20. In this example, $n = 4$, and therefore $n - 1 = 3$, so you divide 20 by 3 to get 6.67. The units here are “points squared,” which obviously makes no sense. Finally, you take the square root of 6.67, to get 2.58. The standard deviation for these four quiz scores is 2.58 points.

Because calculating the standard deviation involves many steps, in most cases you have a computer calculate it for you. However, knowing how to calculate the standard deviation helps you better interpret this statistic and can help you figure out when the statistic may be wrong.



Statisticians divide by $n - 1$ instead of by n in the formula for s so the results have nicer properties that operate on a theoretical plane that's beyond the scope of this book (not the *Twilight Zone* but close; trust me, that's more than you want to know about *that!*).



The standard deviation of an entire population of data is denoted with the Greek letter σ . When I use the term *standard deviation*, I mean s , the sample standard deviation. (When I refer to the population standard deviation, I let you know.)

Interpreting standard deviation

Standard deviation can be difficult to interpret as a single number on its own. Basically, a small standard deviation means that the values in the data set are close to the mean of the data set, on average, and a large standard deviation means that the values in the data set are farther away from the mean, on average.

A small standard deviation can be a goal in certain situations where the results are restricted, for example, in product manufacturing and quality control. A particular type of car part that has to be 2 centimeters in diameter to fit properly had better not have a very big standard deviation during the manufacturing process. A big standard deviation in this case would mean that lots of parts end up in the trash because they don't fit right; either that or the cars will have problems down the road.

But in situations where you just observe and record data, a large standard deviation isn't necessarily a bad thing; it just reflects a large amount of variation in the group that is being studied. For example, if you look at salaries for everyone in a certain company, including everyone from the student intern to the CEO, the standard deviation may be very large. On the other hand, if you narrow the group down by looking only at the student interns, the standard deviation is smaller, because the individuals within this group have salaries that are less variable. The second data set isn't better, it's just less

variable.

Similar to the mean, outliers affect the standard deviation (after all, the formula for standard deviation includes the mean). In the NBA salaries example, the salaries of the L.A. Lakers in the 2009–2010 season (shown in Table 5-2) range from the highest, \$23,034,375 (Kobe Bryant) down to \$959,111 (Didier Ilunga-Mbenga and Josh Powell). Lots of variation, to be sure! The standard deviation of the salaries for this team turns out to be \$6,567,405; it's almost as large as the average. However, as you may guess, if you remove Kobe Bryant's salary from the data set, the standard deviation decreases because the remaining salaries are more concentrated around the mean. The standard deviation becomes \$4,671,508.



Watch for the units when determining whether a standard deviation is large. For example, a standard deviation of 2 in units of years is equivalent to a standard deviation of 24 in units of months. Also look at the value of the mean when putting standard deviation into perspective. If the average number of Internet newsgroups that a user posts to is 5.2 and the standard deviation is 3.4, that's a lot of variation, relatively speaking. But if you're talking about the age of the newsgroup users where the mean is 25.6 years, that same standard deviation of 3.4 would be comparatively smaller.

Understanding properties of standard deviation

Here are some properties that can help you when interpreting a standard deviation:

- ✓ The standard deviation can never be a negative number, due to the way it's calculated and the fact that it measures a distance (distances are never negative numbers).
- ✓ The smallest possible value for the standard deviation is 0, and that happens only in contrived situations where every single number in the data set is exactly the same (no deviation).
- ✓ The standard deviation is affected by outliers (extremely low or extremely high numbers in the data set). That's because the standard deviation is based on the distance from the *mean*. And remember, the mean is also affected by outliers.
- ✓ The standard deviation has the same units as the original data.

Lobbying for standard deviation

The standard deviation is a commonly used statistic, but it doesn't often get the attention it deserves. Although the mean and median are out there in common sight in the everyday media, you rarely see them accompanied by any measure of how diverse

that data set was, and so you are getting only part of the story. In fact, you could be missing the most interesting part of the story.

Without standard deviation, you can't get a handle on whether the data are close to the average (as are the diameters of car parts that come off of a conveyor belt when everything is operating correctly) or whether the data are spread out over a wide range (as are house prices and income levels in the U.S.).

For example if someone told you that the average starting salary for someone working at Company Statistix is \$70,000, you may think, "Wow! That's great." But if the standard deviation for starting salaries at Company Statistix is \$20,000, that's a lot of variation in terms of how much money you can make, so the average starting salary of \$70,000 isn't as informative in the end, is it?

On the other hand, if the standard deviation was only \$5,000, you would have a much better idea of what to expect for a starting salary at that company. Which is more appealing? That's a decision each person has to make; however it'll be a much more informed decision once you realize standard deviation matters.

Without the standard deviation, you can't compare two data sets effectively. Suppose two sets of data have the same average; does that mean that the data sets must be exactly the same? Not at all. For example, the data sets 199, 200, 201; and 0, 200, 400 both have the same average (200) yet they have very different standard deviations. The first data set has a *very* small standard deviation ($s=1$) compared to the second data set ($s=200$).

References to the standard deviation may become more commonplace in the media as more and more people (like you, for example) discover what the standard deviation can tell them about a set of results and start asking for it. In your career, you are likely to see the standard deviation reported and used as well.

Being out of range

The range is another statistic that some folks use to measure diversity in a data set. The *range* is the largest value in the data set minus the smallest value in the data set. It's easy to find; all you do is put the numbers in order (from smallest to largest) and do a quick subtraction. Maybe that's why the range is used so often; it certainly isn't because of its interpretative value.



The range of a data set is almost meaningless. It depends on only two numbers in the data set, both of which may reflect extreme values (outliers). My advice is to ignore the range and find the standard deviation, which is a more informative measure of the variation in the data set because it involves all the values. Or you can also calculate another statistic called the *interquartile range*, which is similar to the range with an important difference — it eliminates outlier and skewness issues by

only looking at the middle 50% of the data and finding the range for those values. The section “Exploring interquartile range” at the end of this chapter gives you more details.

Examining the Empirical Rule (68-95-99.7)

Putting a measure of center (such as the mean or median) together with a measure of variation (such as standard deviation or interquartile range) is a good way to describe the values in a population. In the case where the data are in the shape of a bell curve (that is, they have a normal distribution; see Chapter 9), the population mean and standard deviation are the combination of choice, and a special rule links them together to get some pretty detailed information about the population as a whole.

The *Empirical Rule* says that if a population has a normal distribution with population mean μ and standard deviation σ , then:

- ✓ About 68% of the values lie within 1 standard deviation of the mean (or between the mean minus 1 times the standard deviation, and the mean plus 1 times the standard deviation). In statistical notation, this is represented as $\mu \pm 1\sigma$.
- ✓ About 95% of the values lie within 2 standard deviations of the mean (or between the mean minus 2 times the standard deviation, and the mean plus 2 times the standard deviation). The statistical notation for this is $\mu \pm 2\sigma$.
- ✓ About 99.7% of the values lie within 3 standard deviations of the mean (or between the mean minus 3 times the standard deviation and the mean plus 3 times the standard deviation). Statisticians use the following notation to represent this: $\mu \pm 3\sigma$.

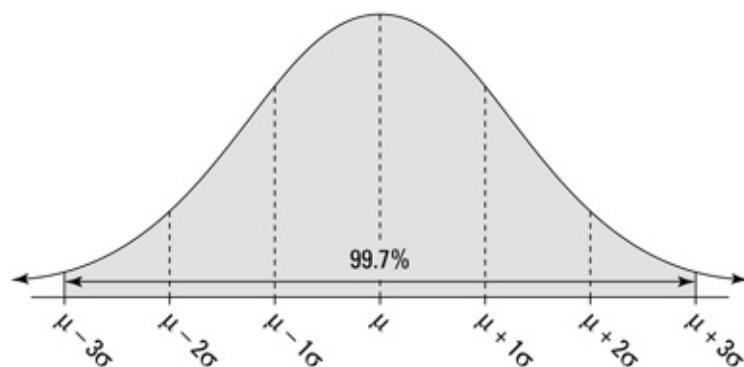
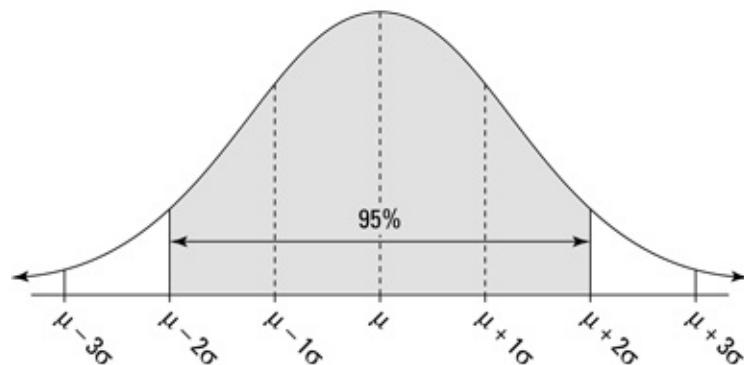
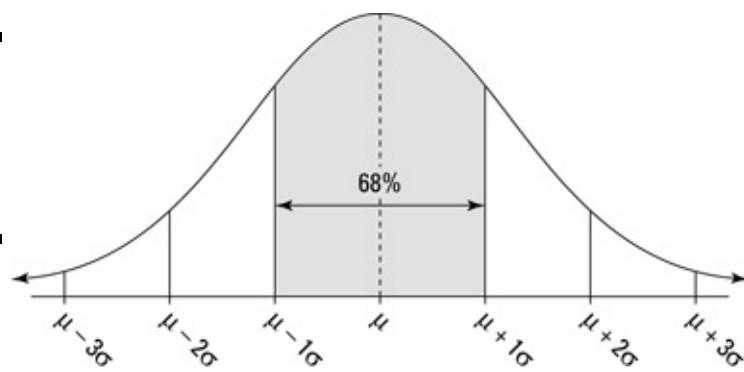


The Empirical Rule is also known as the *68-95-99.7 Rule*, in correspondence with those three properties. It's used to describe a population rather than a sample, but you can also use it to help you decide whether a sample of data came from a normal distribution. If a sample is large enough and you can see that its histogram looks close to a bell-shape, you can check to see whether the data follow the 68-95-99.7 percent specifications. If yes, it's reasonable to conclude the data came from a normal distribution. This is huge because the normal distribution has lots of perks, as you can see in Chapter 9.

Figure 5-2 illustrates all three components of the Empirical Rule.

The reason that so many (about 68%) of the values lie within 1 standard deviation of the mean in the Empirical Rule is because when the data are bell-shaped, the majority of the values are mounded up in the middle, close to the mean (as Figure 5-2 shows).

Figure 5-2:
The Empirical
Rule (68%,
95%, and
99.7%).



Adding another standard deviation on either side of the mean increases the percentage from 68 to 95, which is a big jump and gives a good idea of where “most” of the data are located. Most researchers stay with the 95% range (rather than 99.7%) for reporting their results, because increasing the range to 3 standard deviations on either side of the mean (rather than just 2) doesn’t seem worthwhile, just to pick up that last 4.7% of the values.



The Empirical Rule tells you about what percentage of values are within a certain range of the mean, and I need to stress the word *about*. These results are approximations only, and they only apply if the data follow a normal distribution. However, the Empirical Rule is an important result in statistics because the concept of “going out about two standard deviations to get about 95% of the values” is one that you see mentioned often with confidence intervals and hypothesis tests (see Chapters 13 and 14).

Here's an example of using the Empirical Rule to better describe a population whose values have a normal distribution: In a study of how people make friends in cyberspace using newsgroups, the age of the users of an Internet newsgroup was reported to have a mean of 31.65 years, with a standard deviation of 8.61 years. Suppose the data were graphed using a histogram and were found to have a bell-shaped curve similar to what's shown in Figure 5-2.

According to the Empirical Rule, about 68% of the newsgroup users had ages within 1 standard deviation (8.61 years) of the mean (31.65 years). So about 68% of the users were between ages $31.65 - 8.61$ years and $31.65 + 8.61$ years, or between 23.04 and 40.26 years. About 95% of the newsgroup users were between the ages of $31.65 - 2(8.61)$, and $31.65 + 2(8.61)$, or between 14.43 and 48.87 years. Finally, about 99.7% of the newsgroup users' ages were between $31.65 - 3(8.61)$ and $31.65 + 3(8.61)$, or between 5.82 and 57.48 years.

This application of the rule gives you a much better idea about what's happening in this data set than just looking at the mean, doesn't it? As you can see, the mean and standard deviation used together add value to your results; plugging these values into the Empirical Rule allows you to report ranges for "most" of the data yourself.



Remember, the condition for being able to use the Empirical Rule is that the data have a normal distribution. If that's not the case (or if you don't know what the shape actually is), you can't use it. To describe your data in these cases, you can use percentiles, which represent certain cutoff points in the data (see the later section "Gathering a five-number summary").

Measuring Relative Standing with Percentiles

Sometimes the precise values of the mean, median, and standard deviation just don't matter, and all you are interested in is where you stand compared to the rest of the herd. In this situation, you need a statistic that reports *relative standing*, and that statistic is called a percentile. The k^{th} percentile is a number in the data set that splits the data into two pieces: The lower piece contains k percent of the data, and the upper piece contains the rest of the data (which amounts to $[100 - k]$ percent, because the total amount of data is 100%). **Note:** k is any number between 1 and 100.



The median is the 50th percentile: The point in the data where 50% of the data fall below that point, and 50% fall above it.

In this section, you find out how to calculate, interpret, and put together percentiles to help you uncover the story behind a data set.

Calculating percentiles

To calculate the k^{th} percentile (where k is any number between one and one hundred), do the following steps:

- 1. Order all the numbers in the data set from smallest to largest.**
- 2. Multiply k percent times the total number of numbers, n .**
- 3a. If your result from Step 2 is a whole number, go to Step 4. If the result from Step 2 is not a whole number, round it up to the nearest whole number and go to Step 3b.**
- 3b. Count the numbers in your data set from left to right (from the smallest to the largest number) until you reach the value indicated by Step 3a. The corresponding value in your data set is the k^{th} percentile.**
- 4. Count the numbers in your data set from left to right until you reach the one indicated by Step 2. The k^{th} percentile is the average of that corresponding value in your data set and the value that directly follows it.**

For example, suppose you have 25 test scores, and in order from lowest to highest they look like this: 43, 54, 56, 61, 62, 66, 68, 69, 69, 70, 71, 72, 77, 78, 79, 85, 87, 88, 89, 93, 95, 96, 98, 99, 99. To find the 90th percentile for these (ordered) scores, start by multiplying 90% times the total number of scores, which gives $90\% * 25 = 0.90 * 25 = 22.5$. Rounding up to the nearest whole number, you get 23.

Counting from left to right (from the smallest to the largest number in the data set), you go until you find the 23rd number in the data set. That number is 98, and it's the 90th percentile for this data set.

Now say you want to find the 20th percentile. Start by taking $0.20 * 25 = 5$; this is a whole number, so proceed from Step 3a to Step 4, which tells us the 20th percentile is the average of the 5th and 6th numbers in the ordered data set (62 and 66). The 20th percentile then comes to $(62 + 66) \div 2 = 64$. The median (the 50th percentile) for the test scores is the 13th score: 77.



There is no single definitive formula for calculating percentiles. The formula here is designed to make finding the percentile easier and more intuitive, especially if you're doing the work by hand; however, other formulas are used when you're working with technology. The results you get using various methods may differ but not by much.

Interpreting percentiles

Percentiles report the relative standing of a particular value within a data set. If that's

what you're most interested in, the actual mean and standard deviation of the data set are not important, and neither is the actual data value. What's important is where you stand — not in relation to the mean, but in relation to everyone else: That's what a percentile gives you.

For example, in the case of exam scores, who cares what the mean is, as long as you scored better than most of the class? Who knows, it may have been an impossible exam and 40 points out of 100 was a great score (that happened to me in an advanced math class once; heaven forbid this should ever happen to you!). In this case, your score itself is meaningless, but your percentile tells you everything.

Suppose your exam score is better than 90% of the rest of the class. That means your exam score is at the 90th percentile (so $k = 90$), which hopefully gets you an A. Conversely, if your score is at the 10th percentile (which would never happen to you, because you're such an excellent student), then $k = 10$; that means only 10% of the other scores are below yours, and 90% of them are above yours; in this case an A is not in your future.

A nice property of percentiles is they have a universal interpretation: Being at the 95th percentile means the same thing no matter if you are looking at exam scores or weights of packages sent through the postal service; the 95th percentile always means 95% of the other values lie below yours, and 5% lie above it. This also allows you to fairly compare two data sets that have different means and standard deviations (like ACT scores in reading versus math). It evens the playing field and gives you a way to compare apples to oranges, so to speak.



A percentile is *not* a percent; a percentile is a number (or the average of two numbers) in the data set that marks a certain percentage of the way through the data. Suppose your score on the GRE was reported to be the 80th percentile. This doesn't mean you scored 80% of the questions correctly. It means that 80% of the students' scores were lower than yours and 20% of the students' scores were higher than yours.



A high percentile doesn't always constitute a good thing. For example, if your city is at the 90th percentile in terms of crime rate compared to cities of the same size, that means 90% of cities similar to yours have a crime rate that is lower than yours, which is not good for you. Another example is golf scores; a low score in golf is a good thing, so being at the 80th percentile with your score wouldn't qualify you for the PGA tour, let's just say that.

Comparing household incomes

The U.S. government often reports percentiles among its data summaries. For example,

the U.S. Census Bureau reported the median (the 50th percentile) household income for 2001 to be \$42,228, and in 2007 it was reported to be \$50,233. The Bureau also reports various percentiles for household income for each year, including the 10th, 20th, 50th, 80th, 90th, and 95th. Table 5-3 shows the values of each of these percentiles for both 2001 and 2007.

Table 5-3

U.S. Household Income (2001 versus 2007)

Percentile	2001 Household Income	2007 Household Income
10th	\$10,913	\$12,162
20th	\$17,970	\$20,291
50th	\$42,228	\$50,233
80th	\$83,500	\$100,000
90th	\$116,105	\$136,000
95th	\$150,499	\$177,000

Looking at the percentiles for 2001 in Table 5-3, you can see that the bottom half of the incomes are closer together than the top half of the incomes are. The difference between the 20th percentile and the 50th percentile is about \$24,000, whereas the spread between the 50th percentile and the 80th percentile is more like \$41,000. The difference between the 10th and 50th percentiles is only about \$31,000, whereas the difference between the 50th and the 90th percentiles is a whopping \$74,000.

The percentiles for 2007 are all higher than the percentiles for 2001 (which is a good thing!). They are also more spread out. For 2007, the difference between the 20th and 50th percentiles is around \$30,000, and from the 50th to the 80th it's approximately \$50,000; both of these differences are larger than for 2001. Similarly, the 10th percentile is farther from the 50th (about \$38,000 difference) in 2007 compared to 2001, and the 50th is farther from the 90th (by about \$86,000) in 2007, compared to 2001. These results tell us that incomes are increasing in general at all levels between 2001 and 2007, but the gap is widening between those levels. For example, the 10th percentile for income in 2001 was \$10,913 (as seen in Table 5-3), compared to \$12,162 in 2007; this represents about an 11 percent increase (subtract the two and divide by 10,913). Now compare the 95th percentiles for 2007 versus 2001; the increase is almost 18%. Now, technically, you may want to adjust the 2001 values for inflation, but you get the basic idea.



Percentage changes affect the variability in a data set. For example, when salary raises are given on a percentage basis, the diversity in the salaries also increases; it's the "rich get richer" idea. The guy making \$30,000 gets a 10 percent raise and his salary goes up to \$33,000 (an increase of \$3,000); but the guy making \$300,000 gets a 10 percent raise and now makes \$330,000 (a difference of \$30,000). So when you first get hired for a new job, negotiate the highest possible salary you can because your raises that follow will also net a higher amount.

Examining ACT Scores

Each year millions of U.S. high school students take a nationally administered ACT exam as part of the process of applying for colleges. The test is designed to assess college readiness in the areas of English, Math, Reading, and Science. Each test has a possible score of 36 points.

ACT does not release the average or standard deviation of the test scores for a given exam. (That would be a real hassle if they did, because these statistics can change from exam to exam, and people would complain that this exam was harder than that exam when the actual scores are not relevant.) To avoid these issues, and for other reasons, ACT reports test results using percentiles.

Percentiles are usually reported in the form of a predetermined list. For example, the U.S. Census Bureau reports the 10th, 20th, 50th, 80th, 90th, and 95th percentiles for household income (as shown in Table 5-3). However, ACT uses percentiles in a different way. Rather than reporting the exam scores corresponding to a premade list of percentiles, they list each possible exam score and report its corresponding percentile, whatever that turns out to be. That way, to find out where you stand, you just look up your score and you'll find out your percentile.

Table 5-4 shows the 2009 percentiles for the scores on the Mathematics and Reading ACT exams. To interpret an exam score, find the row corresponding to the score and the column for the exam area (for example, Reading). Intersect row and column and you find out which percentile your score represents; in other words, you see what percentage of your fellow exam-taking comrades scored lower than you.

Table 5-4**Percentiles for All Possible ACT Exam Scores in Math and Reading**

ACT Score	Mathematics Percentile	Reading Percentile
34–36	99	99
33	98	97
32	97	95
31	96	93
30	95	91
29	93	88
28	91	85
27	88	81
26	84	78
25	79	74
24	74	70
23	68	65
22	62	59
21	57	54
20	52	47
19	47	41
18	40	34
17	33	30
16	24	24
15	14	19
14	06	14
13	02	09
12	01	06
11	01	03
1–10	01	01

For example, suppose you scored 30 on the Math exam; in Table 5-4 you look at the row for 30 in the column for Math; you see your score is at the 95th percentile. In other words 95% of the students scored lower than you, and only 5% scored higher than you.

Now suppose you also scored a 30 on the Reading exam. Just because a score of 30 represents the 95th percentile for Math doesn't necessarily mean a score of 30 is at the 95th percentile for Reading as well. (It's probably reasonable to expect that fewer people score 30 or higher on the Math exam than on the Reading exam.)

To test my theory, look at column 3 of Table 5-4 in the row for a score of 30. You see that a score of 30 on the Reading exam puts you at the 91st percentile — not quite as great as your position on the Math exam, but certainly not a bad score.

Gathering a five-number summary

Beyond reporting a single measure of center and/or a single measure of spread, you can create a group of statistics and put them together to get a more detailed description of a data set. The Empirical Rule (as seen in “Examining the Empirical Rule (68-95-99.7)” earlier in this chapter) uses the mean and standard deviation in tandem to describe a bell-shaped data set. In the case where your data are not bell-shaped, you use a different set of statistics (based on percentiles) to describe the big picture of your data. This method involves cutting the data into four pieces (with an equal amount of data in each piece) and reporting the resulting five cutoff points that separate these pieces. These cutoff points are represented by a set of five statistics that describe how the data are laid out.

The *five-number summary* is a set of five descriptive statistics that divide the data set into four equal sections. The five numbers in a five-number summary are:

1. The *minimum* (smallest) number in the data set
2. The *25th percentile* (also known as *the first quartile*, or Q_1)
3. The *median* (50th percentile)
4. The *75th percentile* (also known as *the third quartile*, or Q_3)
5. The *maximum* (largest) number in the data set

For example, suppose you want to find the five-number summary of the following 25 (ordered) exam scores: 43, 54, 56, 61, 62, 66, 68, 69, 69, 70, 71, 72, 77, 78, 79, 85, 87, 88, 89, 93, 95, 96, 98, 99, 99. The minimum is 43, the maximum is 99, and the median is the number directly in the middle, 77.

To find Q_1 and Q_3 you use the steps shown in the section “Calculating percentiles,” with $n = 25$. Step 1 is done because the data are ordered. For Step 2, since Q_1 is the 25th percentile, multiply $0.25 * 25 = 6.25$. This is not a whole number, so Step 3a says to round it up to 7 and proceed to Step 3b.

Following Step 3b, you count from left to right in the data set until you reach the 7th number, 68; this is Q_1 . For Q_3 (the 75th percentile) you multiply $0.75 * 25 = 18.75$, which you round up to 19. The 19th number on the list is 89, so that’s Q_3 . Putting it all together, the five-number summary for these 25 test scores is 43, 68, 77, 89, and 99. To best interpret a five-number summary, you can use a boxplot; see Chapter 7 for details.

Exploring interquartile range

The purpose of the five-number summary is to give descriptive statistics for center, variation, and relative standing all in one shot. The measure of center in the five-number summary is the median, and the first quartile, median, and third quartiles are measures

of relative standing.

To obtain a measure of variation based on the five-number summary, you can find what's called the *interquartile range* (or *IQR*). The *IQR* equals $Q_3 - Q_1$ (that is, the 75th percentile minus the 25th percentile) and reflects the distance taken up by the innermost 50% of the data. If the *IQR* is small, you know a lot of data are close to the median. If the *IQR* is large, you know the data are more spread out from the median. The *IQR* for the test scores data set is $89 - 68 = 21$, which is fairly large, seeing as how test scores only go from 0 to 100.



The interquartile range is a much better measure of variation than the regular range (maximum value minus minimum value; see the section “Being out of range” earlier in this chapter). That’s because the interquartile range doesn’t take outliers into account; it cuts them out of the data set by only focusing on the distance within the middle 50 percent of the data (that is, between the 25th and 75th percentiles).



Descriptive statistics that are well chosen and used correctly can tell you a great deal about a data set, such as where the center is located, how diverse the data are, and where a good portion of the data lies. However, descriptive statistics can't tell you everything about the data, and in some cases they can be misleading. Be on the lookout for situations where a different statistic would be more appropriate (for example, the median describes center more fairly than the mean when the data is skewed), and keep your eyes peeled for situations where critical statistics are missing (for example, when a mean is reported without a corresponding standard deviation).

Chapter 6

Getting the Picture: Graphing Categorical Data

In This Chapter

- ▶ Making data displays for categorical data
- ▶ Interpreting and critiquing charts and graphs

Data displays, especially charts and graphs, seem to be everywhere, showing everything from election results, broken down by every conceivable characteristic, to how the stock market has fared over the past few years (months, weeks, days, minutes). We're living in an instant gratification, fast-information society; everyone wants to know the bottom line and be spared the details.

The abundance of graphs and charts is not necessarily a bad thing, but you have to be careful; some of them are incorrect or even misleading (sometimes intentionally and sometimes by accident), and you have to know what to look for.

This chapter is about graphs involving *categorical data* (data that places individuals into groups or categories, such as gender, opinion, or whether a patient takes medication every day. Here you find out how to read and make sense of these data displays and get some tips for evaluating them and spotting problems. (**Note:** Data displays for *numerical data*, such as weight, exam score, or the *number* of pills taken by a patient each day, come in Chapter 7.)

The most common types of data displays for categorical data are pie charts and bar graphs. In this chapter, I present examples of each type of data display and share some thoughts on interpretation and tips for critically evaluating each type.

Take Another Little Piece of My Pie Chart

A pie chart takes categorical data and breaks them down by group, showing the percentage of individuals that fall into each group. Because a pie chart takes on the shape of a circle, the “slices” that represent each group can easily be compared and contrasted.



Because each individual in the study falls into one and only one category, the sum of all the slices of the pie should be 100% or close to it (subject to a bit of rounding off). However, just in case, keep your eyes open for pie charts whose percentages just don't add up.

Tallying personal expenses

When you spend your money, what do you spend it on? What are your top three expenses? According to the U.S. Bureau of Labor Statistics 2008 Consumer Expenditure Survey, the top six sources of consumer expenditures in the U.S. were housing (33.9%), transportation (17.0%), food (12.8%), personal insurance and pensions (11.1%), healthcare (5.9%), and entertainment (5.6%). These six categories make up over 85% of average consumer expenses. (Although the exact percentages change from year to year, the list of the top six items remains the same.)

Figure 6-1 summarizes the 2008 U.S. expenditures in a pie chart. Notice that the “Other” category is a bit large in this chart (13.7%). However, with so many other possible expenditures out there (including this book), each one would only get a tiny slice of the pie for itself, and the resulting pie chart would be a mess. In this case, it is too difficult to break “Other” down further. (But in many other cases you can.)



Ideally, a pie chart shouldn't have too many slices because a large number of slices distracts the reader from the main point(s) the pie chart is trying to relay. However, lumping the remaining categories into one slice that's one of the largest in the whole pie chart leaves readers wondering what's included in that particular slice. With charts and graphs, doing it right is a delicate balance.

Bringing in a lotto revenue

State lotteries bring in a great deal of revenue, and they also return a large portion of the money received, with some of the revenues going to prizes and some being allocated to state programs such as education. Where does lottery revenue come from? Figure 6-2 is a pie chart showing the types of games and their percentage of revenue as recently reported by Ohio's state lottery. (Note the slices don't sum to 100% exactly due to slight rounding error.)

You can see by the pie chart in Figure 6-2 that 49.3% of the lottery sales revenue comes from the instant (scratch-off) games. The rest come from various lottery-type games in which players choose a set of numbers and win if a certain number of their numbers match those chosen by the lottery.

Figure 6-1:
Pie chart
showing how
people in the
U.S. spend
their money.

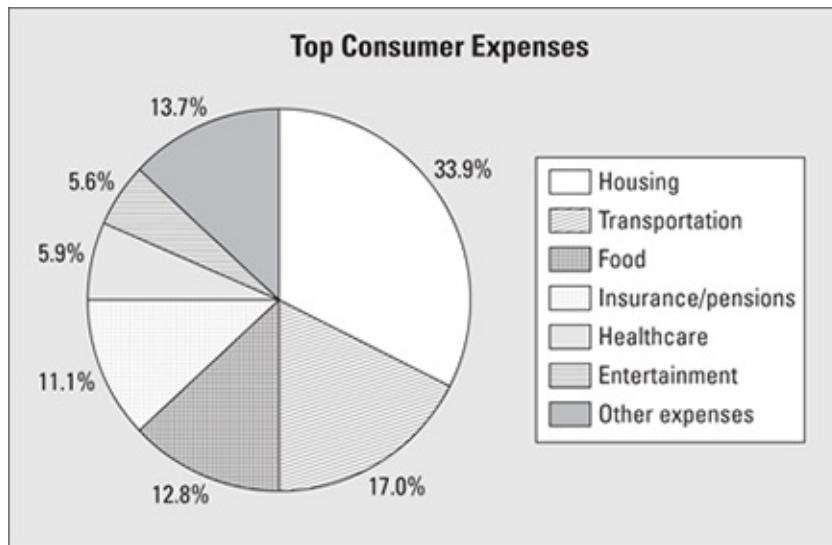
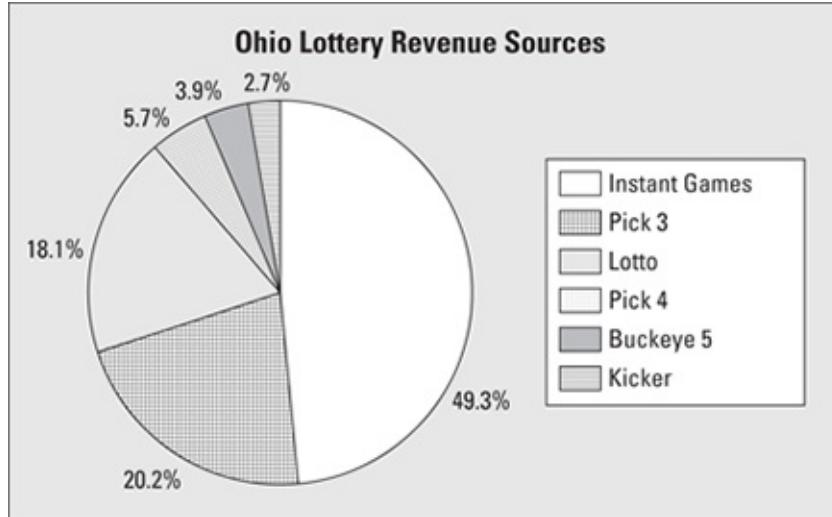


Figure 6-2:
Pie chart
breaking
down a
state's lottery
revenue.



Notice that this pie chart doesn't tell you *how much* money came in, only *what percentage* of the money came from each type of game. About half the money (49.3%) came from instant scratch-off games; does this revenue represent a million dollars, two million dollars, ten million dollars, or more? You can't answer these questions without knowing the total amount of revenue dollars.

I was, however, able to find this information on another chart provided by the lottery Web site: The total revenue (over a 10-year period) was reported as “1,983.1 million dollars” — which you also know as 1.9831 billion dollars. Because 49.3% of sales came from instant games, they therefore represent sales revenue of \$977,668,300 over a 10-year period. That's a lot of (or dare I say a “lotto”) scratching.

Ordering takeout

It's also important to watch for totals when examining a pie chart from a survey. A newspaper I read reported the latest results of a “people poll.” They asked, “What is your favorite night to order takeout for dinner?” The results are shown in a pie chart (see Figure 6-3).

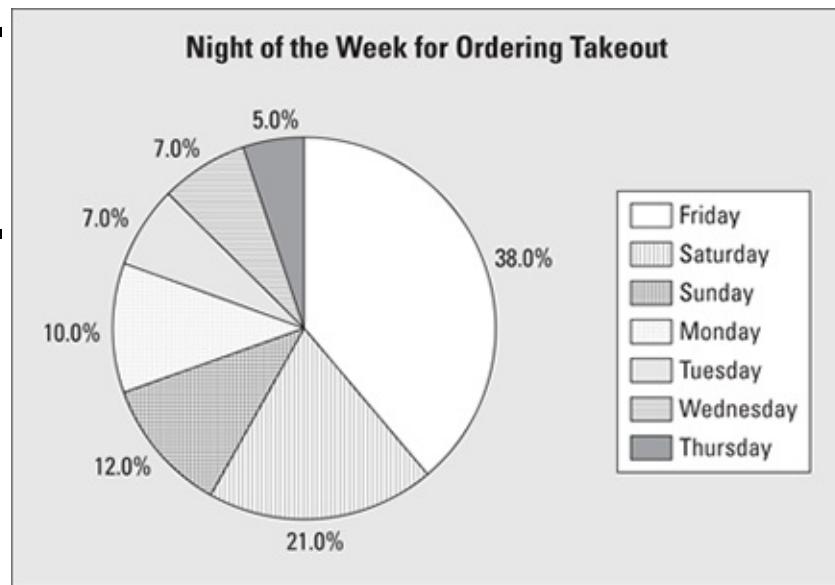
You can clearly see that Friday night is the most popular night for ordering takeout (and that result makes sense) with decreasing demand moving from Saturday through Monday. The actual percentages shown in Figure 6-3 really only apply to the people who were surveyed; how close these results mimic the population depends on many factors, one of which is sample size. But unfortunately, sample size is not included as part of this graph. (For example, it would be nice to see “ $n = XXX$ ” below the title; where n represents sample size.)

Without knowing the sample size, you can't tell how accurate the information is. Which results would you find to be more accurate — those based on 25 people, 250 people, or 2,500 people? When you see the number 10%, you don't know if it's 10 out of 100, 100 out of 1,000, or even 1 out of 10. To statisticians, $1 \div 10$ is not the same as $100 \div 1,000$, even though they both represent 10%. (Don't tell that to mathematicians — they'll think you're nuts!)



Pie charts often don't include mention of the total sample size. Always check for the sample size, especially if the results are very important to you; don't assume it's large! If you don't see the sample size, go to the source of the data and ask for it.

Figure 6-3:
Pie chart for
takeout food
survey
results.



Projecting age trends

The U.S. Census Bureau provides an almost unlimited amount of data, statistics, and graphics about the U.S. population, including the past, present, and projections for the future. It often makes comparisons between years in order to look for changes and trends.

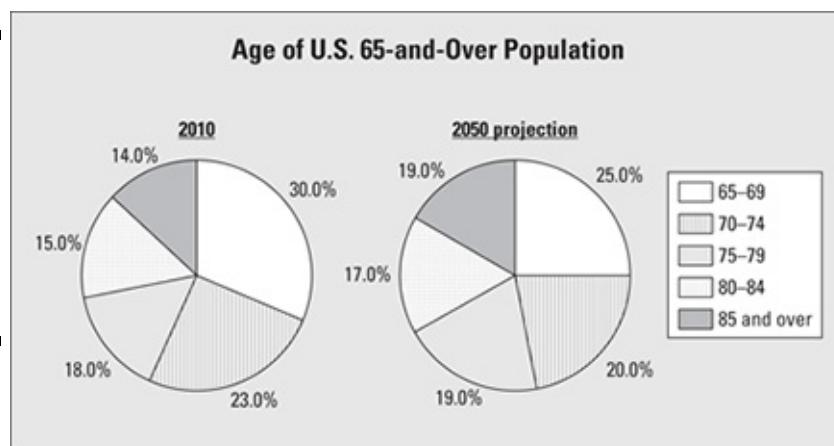
One recent Census Bureau population report looked at what it calls the “older U.S. population” (by the government’s definition, this means people 65 years old or over). Age was broken into the following groups: 65–69 years, 70–74 years, 75–79 years, 80–84

years, and 85 and over. The Bureau calculated and reported the percentage in each age group for the year 2010 and made projections for the percentage in each age group for the year 2050.

I made side-by-side pie charts for the years 2010 versus 2050 (projections) to make comparisons; you can see the results in Figure 6-4. The percentage of the older population in each age group for 2010 is shown in one pie chart, and alongside it is a pie chart of the projected percentage for each age group for 2050 (based on the current age of the entire U.S. population, birth and death rates, and other variables).

If you compare the sizes of the slices from one graph to the other in Figure 6-4, you see that the slices for corresponding age groups are larger for the 2050 projections (compared to 2010) as the age groups get older, and the slices are smaller for the 2050 projections (compared to 2010) as the age groups get younger. For example the 65–69 age group decreases from 30% in 2010 to a projected 25% in 2050; while the 85-and-over age group increases from 14% in 2010 to 19% projected for 2050.

Figure 6-4:
Side-by-side
pie charts on
the aging
population,
2010 versus
2050
projections.



The results from Figure 6-4 indicate a shift in the ages of the population toward the older categories. From there, the medical and social research communities can examine the ramifications of this trend in terms of healthcare, assisted living, social security, and so on.



The operative words here are *if the trend continues*. As you know, many variables affect population size, and you need to take those into account when interpreting these projections into the future. The U.S. government always points out caveats like this in their reports; it is very diligent about that.



The pie charts in Figure 6-4 work well for comparing groups because they are side-by-side on the same graph, using the same coding for the age groups in each, and their slices are in the same order for both as you move clockwise around the graphs. They aren't all scrambled up on each graph so you have to hunt for a certain age group on each graph separately.

Evaluating a pie chart

The following tips help you taste test a pie chart for statistical correctness:

- ✓ Check to be sure the percentages add up to 100% or very close to it (any round-off error should be very small).
- ✓ Beware of slices of the pie called “Other” that are larger than many of the other slices.
- ✓ Look for a reported total number of units (people, dollar amounts, and so on) so that you can determine (in essence) how “big” the pie was before being divided up into the slices that you’re looking at.
- ✓ Avoid three-dimensional pie charts; they don’t show the slices in their proper proportions. The slices in front look larger than they should.

Raising the Bar on Bar Graphs

A *bar graph* (or *bar chart*) is perhaps the most common data display used by the media. Like a pie chart, a bar graph breaks categorical data down by group. Unlike a pie chart, it represents these amounts by using bars of different lengths; whereas a pie chart most often reports the amount in each group as percentages, a bar graph uses either the number of individuals in each group (also called the *frequency*) or the percentage in each group (called the *relative frequency*).

Tracking transportation expenses

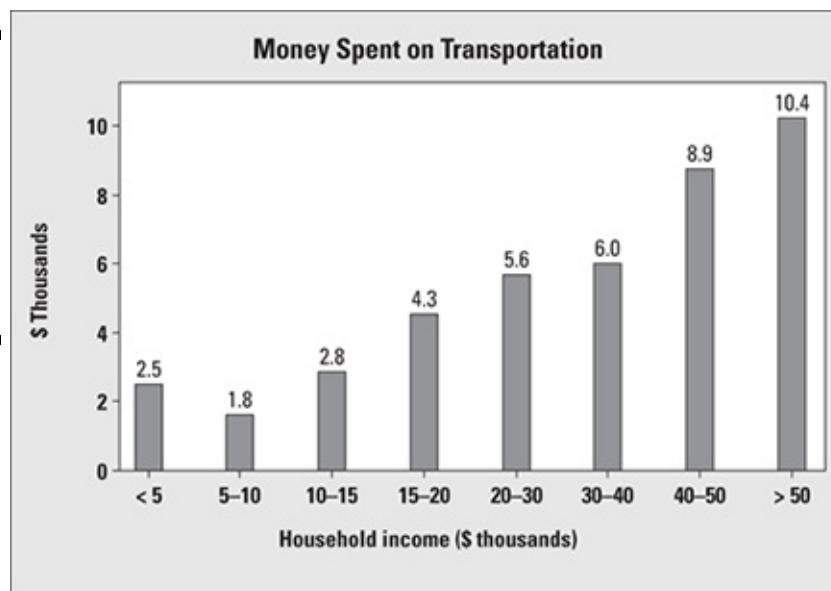
How much of their income do people in the United States spend on transportation to get back and forth to work? It depends on how much money they make. The Bureau of Transportation Statistics (did you know such a department existed?) conducted a study on transportation in the U.S. recently, and many of its findings are presented as bar graphs like the one shown in Figure 6-5.

This particular bar graph shows how much money is spent on transportation for people in different household-income groups. It appears that as household income increases, the total expenditures on transportation also increase. This makes sense, because the more money people have, the more they have available to spend.

But would the bar graph change if you looked at transportation expenditures not in terms of total dollar amounts, but as the percentage of household income? The households in the first group make less than \$5,000 a year and have to spend \$2,500 of it

on transportation. (**Note:** The label reads “2.5,” but because the units are in thousands of dollars, the 2.5 translates into \$2,500.)

Figure 6-5:
Bar graph showing transportation expenses by household income group.



This \$2,500 represents 50% of the annual income of those who make \$5,000 per year; the percentage of the total income is even higher for those who make less than \$5,000 per year. The households earning \$30,000–\$40,000 per year pay \$6,000 per year on transportation, which is between 15% and 20% of their household income. So, although the people making more money spend more dollars on transportation, they don’t spend more as a percentage of their total income. Depending on how you look at expenditures, the bar graph can tell two somewhat different stories.

Another point to check out is the groupings on the graph. The categories for household income as shown aren’t equivalent. For example, each of the first four bars represents household incomes in intervals of \$5,000, but the next three groups increase by \$10,000 each, and the last group contains every household making more than \$50,000 per year. Bar graphs using different-sized intervals to represent numerical values (such as Figure 6-5) make true comparisons between groups more difficult. (However, I’m sure the government has its reasons for reporting the numbers this way; for example, this may be the way income is broken down for tax-related purposes.)

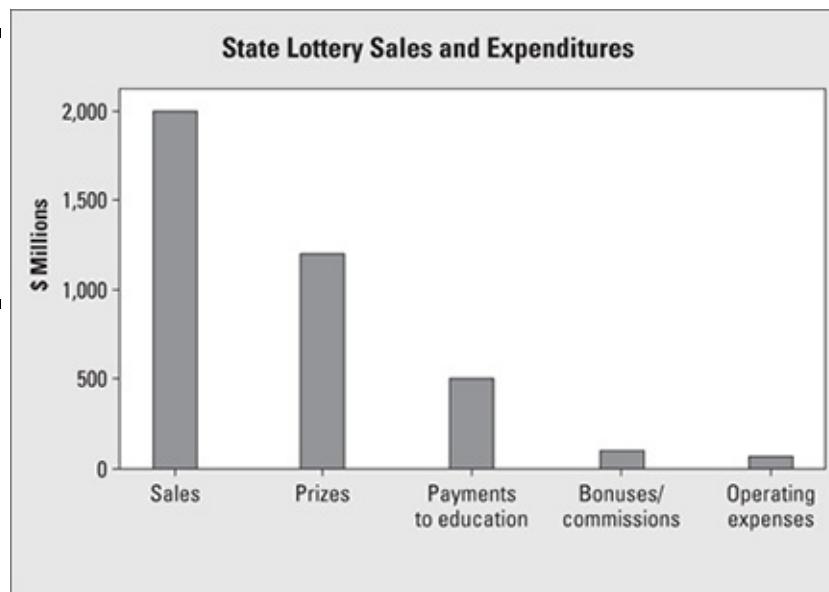
One last thing: Notice that the numerical groupings in Figure 6-5 overlap on the boundaries. For example, \$30,000 appears in both the 5th and 6th bars of the graph. So, if you have a household income of \$30,000, which bar do you fall into? (You can’t tell from Figure 6-5, but I’m sure the instructions are buried in a huge report in the basement of some building in Washington, D.C.) This kind of overlap appears quite frequently in graphs, but you need to know how the borderline values are being treated. For example, the rule may be “Any data lying exactly on a boundary value automatically goes into the bar to its immediate right.” (Looking at Figure 6-5, that puts a household with a \$30,000 income into the 6th bar rather than the 5th.) As long as they are being consistent for each boundary, that’s okay. The alternative, describing the income boundaries for the 5th bar as “\$20,000 to \$29,999.99,” is not an improvement. Along those lines, income data can also

be presented using a histogram (see Chapter 7), which has a slightly different look to it.

Making a lotto profit

That lotteries rake in the bucks is a well-known fact; but they also shell it out. How does it all shake out in terms of profits? Figure 6-6 shows the recent sales and expenditures of a certain state lottery.

Figure 6-6:
Bar graph of
lottery sales
and
expenditures
for a certain
state.



In my opinion, this bar graph needs some additional info from behind the scenes to make it more understandable. The bars in Figure 6-6 don't represent similar types of entities. The first bar represents sales (a form of revenue), and the other bars represent expenditures. The graph would be much clearer if the first bar weren't included; for example, the total sales could be listed as a footnote.

Tipping the scales on a bar graph



Another way a graph can be misleading is through its choice of scale on the frequency/relative frequency axis (that is, the axis where the amounts in each group are reported), and/or its starting value.

By using a “stretched out” scale (for example, having each half inch of a bar represent 10 units versus 50 units), you can stretch the truth, make differences look more dramatic, or exaggerate values. Truth-stretching can also occur if the frequency axis starts out at a number that's very close to where the differences in the heights of the bars start; you are in essence chopping off the bottom of the bars (the less exciting part) and just showing their tops; emphasizing (in a misleading way) where the action is. Not every frequency axis has to start at zero, but watch for situations that elevate the differences.

A good example of a graph with a stretched out scale is seen in Chapter 3, regarding the results of numbers drawn in the “Pick 3” lottery. (You choose three one-digit numbers and if they all match what’s drawn, you win.) In Chapter 3, the percentage of times each number (from 0–9) was drawn is shown in Table 3-2, and the results are displayed in a bar graph in Figure 3-1a. The scale on the graph is stretched and starts at 465, making the differences in the results look larger than they really are; for example, it looks like the number 1 was drawn much less often, whereas the number 2 was drawn much more often, when in reality there is no statistical difference between the percentage of times each number was drawn. (I checked.)

Why was the graph in Figure 3-1a made this way? It might lead people to think they’ve got an inside edge if they choose the number 2 because it’s “on a hot streak”; or they might be led to choose the number 1 because it’s “due to come up.” Both of these theories are wrong, by the way; because the numbers are chosen at random, what happened in the past doesn’t matter. In Figure 3-1b you see a graph that’s been made correctly. (For more examples of where our intuition can go wrong with probability and what the scoop really is, see another of my books, *Probability For Dummies*, also published by Wiley.)

Alternatively, by using a “squeezed down” scale (for example, having each half inch of a bar represent 50 units versus 10 units), you can downplay differences, making results look less dramatic than they actually are. For example, maybe a politician doesn’t want to draw attention to a big increase in crime from the beginning to the end of her term, so she may have the number of crimes of each type shown where each half inch of a bar represents 500 crimes, versus 100 crimes. This squeezes the numbers together and makes differences less noticeable. Her opponent in the next election would go the other way and use a stretched-out scale to emphasize a crime increase in dramatic fashion, and voilà! (Now you know the answer to the question “How can two people talk about the same data and get two different conclusions?” Welcome to the world of politics.)

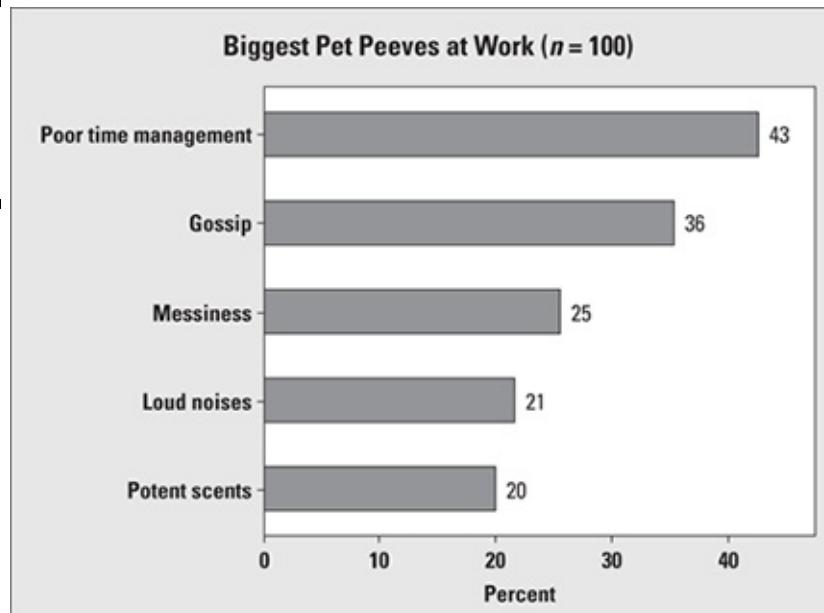


With a pie chart, however, the scale can’t be changed to over-emphasize (or downplay) the results. No matter how you slice up a pie chart, you’re always slicing up a circle, and the proportion of the total pie belonging to any given slice won’t change, even if you make the pie bigger or smaller.

Pondering pet peeves

A recent survey of 100 people with office jobs asked them to report their biggest pet peeves in the workplace. (Before going on, you may want to jot down a couple of yours, just for fun.) A bar graph of the results of the survey is shown in Figure 6-7. Poor time management looks to be the number-one issue for these workers (I hope they didn’t do this survey on company time).

Figure 6-7:
Bar graph for
survey data
with multiple
responses.



Evaluating a bar graph

To raise the statistical bar on bar graphs, check out these tips:

- ✓ Bars that divide up values of a numerical variable (such as income) should be equal in width (if possible) for fair comparison.
- ✓ Be aware of the scale of the bar graph and determine whether it's an appropriate representation of the information.
- ✓ Some bar graphs don't sum to one because they are showing the results of more than one variable; make sure it's clear what's being summarized.
- ✓ Check whether the results are shown as the percentage within each group (relative frequencies) or the number in each group (frequencies).
- ✓ If you see relative frequencies, check for the total sample size — it matters. If you see frequencies, divide each one by the total sample size to get percentages, which are easier to compare.



If you take a look at the percentages shown for each pet peeve listed, you see they don't sum to one. That tells you that each person surveyed was allowed to choose more than one pet peeve (like that would be hard to do); perhaps they were asked to name their top three pet peeves, for example. For this data set and others like it that allow for multiple responses, a pie chart wouldn't be possible (unless you made one for every single pet peeve on the list).

Note that Figure 6-7 is a *horizontal bar graph* (its bars go side to side) as opposed to a *vertical bar graph* (in which bars go up and down, as in Figure 6-6). Either orientation is fine; use whichever one you prefer when you make a bar graph. Do, however, make sure

that you label the axes appropriately and include proper units (such as gender, opinion, or day of the week) where appropriate.

Going by the Numbers: Graphing Numerical Data

In This Chapter

- ▶ Making and interpreting histograms and boxplots for numerical data
 - ▶ Examining time charts for numerical data collected over time
 - ▶ Strategies for spotting misleading and incorrect graphs
-

The main purpose of charts and graphs is to summarize data and display the results to make your point clearly, effectively, and correctly. In this chapter, I present data displays used to summarize *numerical* data — data that represent *counts* (such as the number of pills a patient with diabetes takes per day, or the number of accidents at an intersection per year) or *measurements* (the time it takes you to get to work/school each day, or your blood pressure).

You see examples of how to make, interpret, and evaluate the most common data displays for numerical data: time charts, histograms, and boxplots. I also point out many potential problems that can occur in these graphs, including how people often misread what's there. This information will help you develop important detective skills for quickly spotting misleading graphs.

Handling Histograms

A histogram provides a snapshot of all the data broken down into numerically ordered groups, making it a quick way to get the big picture of the data, in particular, its general shape. In this section you find out how to make and interpret histograms, and how to critique them for correctness and fairness.

Making a histogram

A *histogram* is a special graph applied to data broken down into numerically ordered groups; for example, age groups such as 10–20, 21–30, 31–40, and so on. The bars connect to each other in a histogram — as opposed to a bar graph (Chapter 6) for categorical data, where the bars represent categories that don't have a particular order, and are separated. The height of each bar of a histogram represents either the number of individuals (called the *frequency*) in each group or the percentage of individuals (the *relative frequency*) in each group. Each individual in the data set falls into exactly one bar.



You can make a histogram from any numerical data set; however, you can't determine the actual values of the data set from a histogram because all you know is which group each data value falls into.

An award winning example

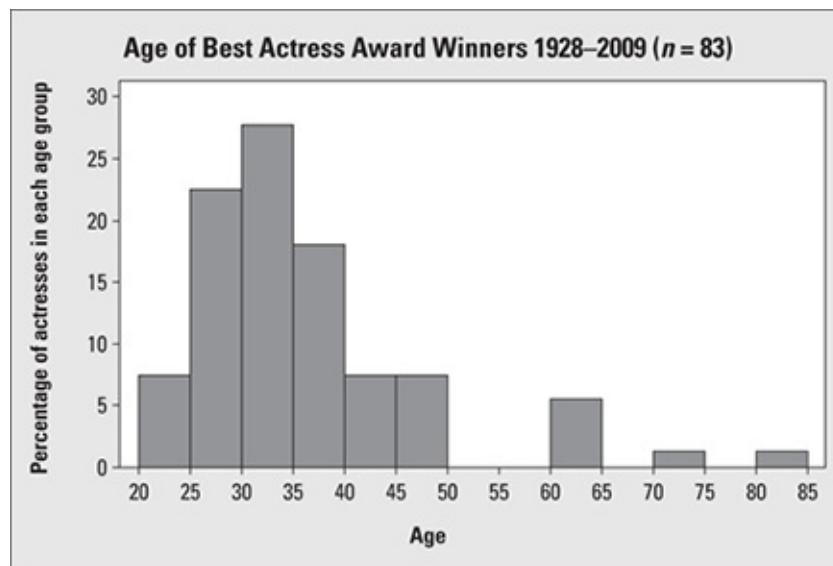
Here's an example of how to create a histogram for all you movie lovers out there (especially those who love old movies). The Academy Awards started in 1928, and one of the most popular categories for this award is Best Actress in a Motion Picture. Table 7-1 shows the winners of the first eight Best Actress Oscars, the years they won (1928–1935), their ages at the time of winning their awards, and the movies they were in. From the table you see the ages range from 22 to 62 — much wider than you may have thought it would be.

Table 7-1 **Ages of Best Actress Oscar Award**
Winners 1928–1935

Year	Winner	Age	Movie
1928	Laura Gainor	22	<i>Sunrise</i>
1929	Mary Pickford	37	<i>Coquette</i>
1930	Norma Shearer	30	<i>The Divorcee</i>
1931	Marie Dressler	62	<i>Min and Bill</i>
1932	Helen Hayes	32	<i>The Sin of Madelon Claudet</i>
1933	Katharine Hepburn	26	<i>Morning Glory</i>
1934	Collette Colbert	31	<i>It Happened One Night</i>
1935	Bette Davis	27	<i>Dangerous</i>

To find out more about the ages of Best Actresses, I expanded my data set to the period 1928–2009. The age variable for this data set is numerical, so you can graph it using a histogram. From there you can answer questions like: What do the ages of these actresses look like? Are they mostly young, old, in between? Are their ages all spread out, or are they similar? Are most of them in a certain age range, with a few outliers (either very young or very old actresses, compared to the others)? To investigate these questions, a histogram of ages of the Best Award actresses is shown in Figure 7-1.

Figure 7-1:
Histogram of
Best Actress
Academy
Award
winners'
ages, 1928–
2009.



Notice that the age groups are shown on the horizontal (*x*) axis. They go by groups of 5 years each: 20–25, 25–30, 30–35, . . . 80–85. The percentage (relative frequency) of actresses in each age group appears on the vertical (*y*) axis. For example, about 27 percent of the actresses were between 30 and 35 years of age when they won their Oscars.

Creating appropriate groups



For Figure 7-1, I used groups of 5 years each in the above example because increments of 5 create natural breaks for years and because it provides enough bars to look for general patterns. You don't have to use this particular grouping, however; you have a bit of poetic license when making a histogram. (However, this freedom allows others to deceive you as you see in the later section "Detecting misleading histograms.") Here are some tips for setting up your histogram:

- ✓ Each data set requires different ranges for its groupings, but you want to avoid ranges that are too wide or too narrow.
 - If a histogram has really wide ranges for its groups, it places all the data into a very small number of bars that make meaningful comparisons impossible.
 - If the histogram has very narrow ranges for its groups, it looks like a big series of tiny bars that cloud the big picture. This can make the data look very choppy with no real pattern.
- ✓ Make sure your groups have equal widths. If one bar is wider than the others, it may contain more data than it should.

One idea that may be appropriate for your histogram is to take the range of the data (largest minus smallest) and divide by 10 to get 10 groupings.

Handling borderline values

In the Academy Award example, what happens if an actress's age lies right on a borderline? For example, in Table 7-1 Norma Shearer was 30 years old in 1930 when she won the Oscar for *The Divorcee*. Does she belong in the 25–30 age group (the lower bar) or the 30–35 age group (the upper bar)?



As long as you are consistent with all the data points, you can either put all the borderline points into their respective lower bars or put all of them into their respective upper bars. The important thing is to pick a direction and be consistent. In Figure 7-1, I went with the convention of putting all borderline values into their respective upper bars — which puts Norma Shearer's age in the 3rd bar, the 30–35 age group of Figure 7-1.

Clarifying the axes

The most complex part of interpreting a histogram for the reader is to get a handle on what's being shown on the x and y axes. Having good descriptive labels on the axes will help. Most statistical software packages label the x -axis using the variable name you provided when you entered your data (for example “age” or “weight”). However, the label for the y -axis isn't as clear. Statistical software packages often label the y -axis of a histogram by writing “frequency” or “percent” by default. These terms can be confusing: frequency or percentage of what?



Clarify the y -axis label on your histogram by changing “frequency” to “number of” and adding the variable name. To modify a label that simply reads “percent,” clarify by writing “percentage of” and the variable. For example, in the histogram of ages of the Best Actress winners shown in Figure 7-1, I labeled the y -axis “Percentage of actresses in each age group.” In the next section you see how to interpret the results from a histogram. How old are those actresses anyway?

Interpreting a histogram



A histogram tells you three main features of numerical data:

- ✓ How the data are distributed among the groups (statisticians call this the *shape* of the data)
- ✓ The amount of variability in the data (statisticians call this the amount of *spread* in the data)
- ✓ Where the center of the data is (statisticians use different measures)

Checking out the shape of the data

One of the features that a histogram can show you is the *shape* of the data — in other words, the manner in which the data fall into the groups. For example, all the data may be exactly the same, in which case the histogram is just one tall bar; or the data might have an equal number in each group; in which case the shape is flat.

Some data sets have a distinct shape. Here are three shapes that stand out:

- ✓ **Symmetric:** A histogram is symmetric if you cut it down the middle and the left-hand and right-hand sides resemble mirror images of each other.

Figure 7-2a shows a symmetric data set; it represents the amount of time each of 50 survey participants took to fill out a certain survey. You see that the histogram is close to symmetric.

- ✓ **Skewed right:** A skewed right histogram looks like a lopsided mound, with a tail going off to the right.

Figure 7-1, showing the ages of the Best Actress Award winners, is skewed right. You see on the right side there are a few actresses whose ages are older than the rest.

- ✓ **Skewed left:** If a histogram is skewed left, it looks like a lopsided mound with a tail going off to the left.

Figure 7-2b shows a histogram of 17 exam scores. The shape is skewed left; you see a few students who scored lower than everyone else.



Following are some particulars about classifying the shape of a data set:

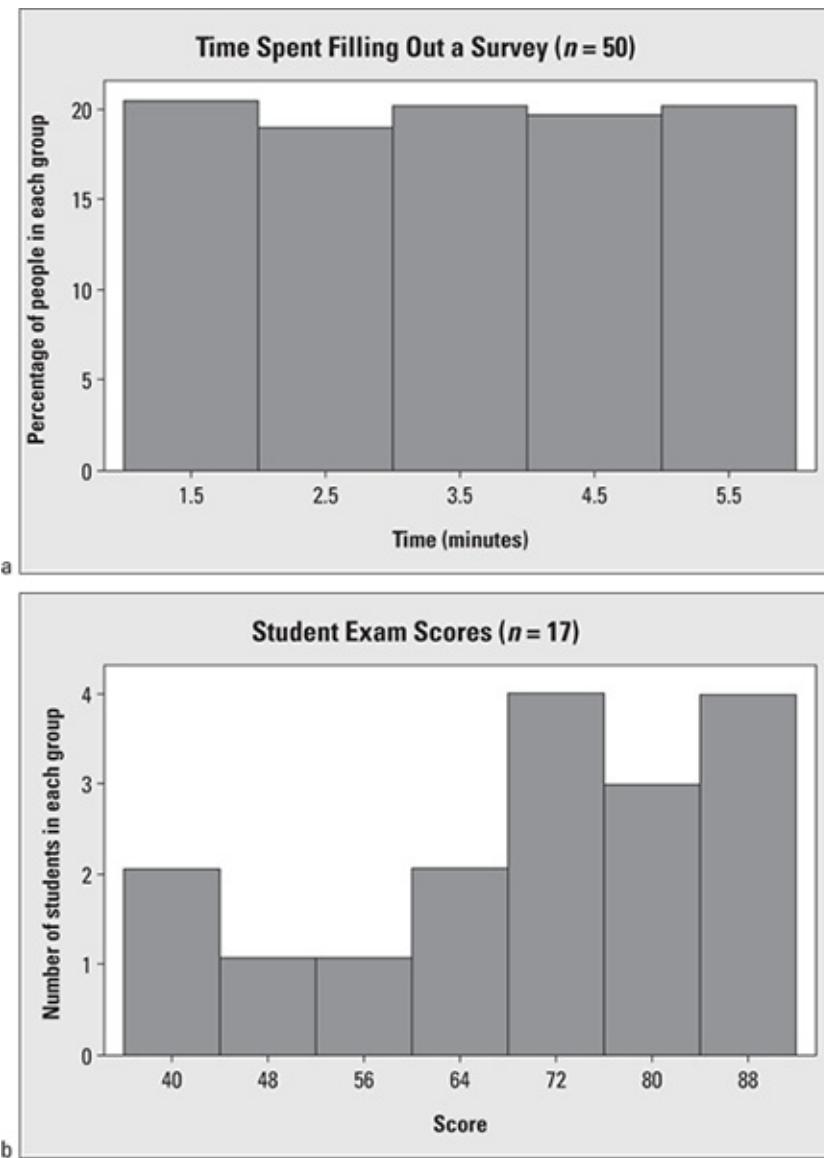
- ✓ **Don't expect symmetric data to have an exact and perfect shape.** Data hardly ever fall into perfect patterns, so you have to decide whether the data shape is close enough to be called symmetric.

If the shape is close enough to symmetric that another person would notice it, and the differences aren't enough to write home about, I'd classify it as symmetric or roughly symmetric. Otherwise, you classify the data as non-symmetric. (More sophisticated statistical procedures exist that actually test data for symmetry, but they're beyond the scope of this book.)

- ✓ **Don't assume that data are skewed if the shape is non-symmetric.** Data sets come in all shapes and sizes, and many of them don't have a distinct shape at all. I include skewness on the list here because it's one of the more common non-symmetric shapes, and it's one of the shapes included in a standard introductory statistics course.

If a data set does turn out to be skewed (or close to it), make sure to denote the direction of the skewness (left or right).

Figure 7-2:
Comparing
the shape of
a) a
symmetric
histogram
and b) a
skewed left
histogram.



As you know from Figure 7-1, the actresses' ages in Figure 7-1 are skewed right. Most of the actresses were between 20 and 50 years of age when they won, with about 27% of them between the ages of 30–35. A few actresses were older when they won their Oscars; about 6 percent were between 60–65 years of age, and less than 4% (total) were 70 years old or over (if you add the percentages from the last two bars in the histogram). The last three bars are what make the data have a shape that is skewed right.

Measuring center: Mean versus median

A histogram gives you a rough idea of where the “center” of the data lies. The word *center* is in quotes because many different statistics are used to designate center. The two most common measures of center are the average (the mean) and the median. (For details on measures of center, see Chapter 5.)



To visualize the average age (the mean), picture the data as people sitting on a teeter-totter. Your objective is to balance it. Because data don't move around, assume the people stay where they are and you move the pivot point (which you can also think of as the hinge or fulcrum) anywhere you want. The mean is the place the pivot point has to be in order to balance the weight on each side of the teeter-totter.

The balancing point of the teeter-totter is affected by the weights of the people on each side, not by the number of people on each side. So the mean is affected by the actual values of the data, rather than the amount of data.

The median is the place where you put the pivot point so you have an equal number of people on each side of the teeter-totter, regardless of their weights. With the same number of people on each side, the teeter-totter wouldn't balance in terms of weight unless the teeter-totter had people with the same total weight on each side. So the median isn't affected by the values of the data, just their location within the data set.



The mean is affected by *outliers*, values in the data set that are away from the rest of the data, on the high end and/or the low end. The median, being the middle number, is not affected by outliers.

Viewing variability: Amount of spread around the mean

You also get a sense of variability in the data by looking at a histogram. For example, if the data are all the same, they are all placed into a single bar, and there is no variability. If an equal amount of data is in each group, the histogram looks flat with the bars close to the same height; this means a fair amount of variability.



The idea of a flat histogram indicating some variability may go against your intuition, and if it does you're not alone. If you're thinking a flat histogram means no variability, you're probably thinking about a time chart, where single numbers are plotted over time (see the section "Tackling Time Charts" later in this chapter). Remember, though, that a histogram doesn't show data over time — it shows all the data at one point in time.

Equally confusing is the idea that a histogram with a big lump in the middle and tails sloping sharply down on each side actually has less variability than a histogram that's straight across. The curves looking like hills in a histogram represent clumps of data that are close together; a flat histogram shows data equally dispersed, with more variability.



Variability in a histogram is higher when the taller bars are more spread out around

the mean and lower when the taller bars are close to the mean.

For the Best Actress Award winners' ages shown in Figure 7-1, you see many actresses are in the age range from 30–35, and most of the ages are between 20–50 years in age, which is quite diverse; then you have those outliers, those few older actresses (I count 7 of them) that spread the data out farther, increasing its overall variability.

The most common statistic used to measure variability in a data set is the *standard deviation*, which in a rough sense measures the average distance that the data lie from the mean. The standard deviation for the Best Actress age data is 11.35 years. (See Chapter 5 for all the details on standard deviation.) A standard deviation of 11.35 years is fairly large in the context of this problem, but the standard deviation is based on average distance from the mean, and the mean is influenced by outliers, so the standard deviation will be as well (see Chapter 5 for more information).

In the later section “Interpreting a boxplot,” I discuss another measure of variability, called the *interquartile range (IQR)*, which is a more appropriate measure of variability when you have skewed data.

Putting numbers with pictures



You can't actually calculate measures of center and variability from the histogram itself because you don't know the exact data values. To add detail to your findings, you should always calculate the basic statistics of center and variation along with your histogram. (All the descriptive statistics you need, and then some, appear in Chapter 5.)

Figure 7-1 is a histogram for the Best Actress ages; you can see it is skewed right. Then for Figure 7-3, I calculated some basic (that is, descriptive) statistics from the data set. Examining these numbers, you find the median age is 33.00 years and the mean age is 35.69 years.

The mean age is higher than the median age because of a few actresses that were quite a bit older than the rest when they won their awards. For example, Jessica Tandy won for her role in *Driving Miss Daisy* when she was 81, and Katharine Hepburn won the Oscar for *On Golden Pond* when she was 74. The relationship between the median and mean confirms the skewness (to the right) found in Figure 7-1.

Figure 7-3:
Descriptive statistics for Best Actress ages (1928–2009).

Descriptive Statistics: Age										
Variable	Total									
		Count	Mean	StDev	Minimum	Q1	Median	Q3	Maximum	IQR
Age	83	35.69	11.35	21.00	28.00	33.00	39.00	81.00	11.00	

Here are some tips for connecting the shape of the histogram (discussed in the previous section) with the mean and median:

✓ **If the histogram is skewed right, the mean is greater than the median.**

This is the case because skewed-right data have a few large values that drive the mean upward but do not affect where the exact middle of the data is (that is, the median). Looking at the histogram of ages of the Best Actress Award winners in Figure 7-1, you see they're skewed right.

✓ **If the histogram is close to symmetric, then the mean and median are close to each other.**



Close to symmetric means it's almost the same on either side; it doesn't need to be exact. *Close* is defined in the context of the data; for example, the numbers 50 and 55 are said to be close if all the values lie between 0 and 1,000, but they are considered to be farther apart if all the values lie between 49 and 56.

The histogram shown in Figure 7-2a is close to symmetric. Its mean and median are both equal to 3.5.

✓ **If the histogram is skewed left, the mean is less than the median.**

This is the case because skewed-left data have a few small values that drive the mean downward but do not affect where the exact middle of the data is (that is, the median).

Figure 7-2b represents the exam scores of 17 students, and the data are skewed left. I calculated the mean and median of the original data set to be 70.41 and 74.00, respectively. The mean is lower than the median due to a few students who scored quite a bit lower than the others. These findings match the general shape of the histogram shown in Figure 7-2b.



The tips for interpreting histograms found in the previous section can also be used the other way around. If for some reason you don't have a histogram of the data, and you only have the mean and median to go by, you compare them to each other to get a rough idea as to the shape of the data set.

✓ If the mean is much larger than the median, the data are generally skewed right; a few values are larger than the rest.

✓ If the mean is much smaller than the median, the data are generally skewed left; a few smaller values bring the mean down.

- ✓ If the mean and median are close, you know the data is fairly balanced, or symmetric, on each side.



Under certain conditions, you can put together the mean and standard deviation to describe a data set in quite a bit of detail. If the data have a normal distribution (a bell-shaped hill in the middle, sloping down at the same rate on each side; see Chapter 5), the Empirical Rule can be applied.

The Empirical Rule (also in Chapter 5) says that if the data have a normal distribution, about 68% of the data lie within 1 standard deviation of the mean, about 95% of the data lie within 2 standard deviations from the mean, and 99.7% of the data lie within 3 standard deviations of the mean. These percentages are custom-made for the normal distribution (bell-shaped data) only and can't be used for data sets of other shapes.

Detecting misleading histograms

There are no hard and fast rules for how to create a histogram; the person making the graph gets to choose the groupings on the x -axis as well as the scale and starting and ending points on the y -axis. Just because there is an element of choice, however, doesn't mean every choice is appropriate; in fact, a histogram can be made to be misleading in many ways. In the following sections, you see examples of misleading histograms and how to spot them.

Missing the mark with too few groups

Although the number of groups you use for a histogram is up to the discretion of the person making the graph, there is such a thing as going overboard, either by having way too few bars, with everything lumped together, or by having way too many bars, where every little difference is magnified.



To decide how many bars a histogram should have, I take a good look at the groupings used to form the bars on the x -axis and see if they make sense. For example, it doesn't make sense to talk about exam scores in groups of 2 points; that's too much detail — too many bars. On the other hand, it doesn't make sense to group actresses' ages by intervals of 20 years; that's not descriptive enough.

Figures 7-4 and 7-5 illustrate this point. Each histogram summarizes $n = 222$ observations of the amount of time between eruptions of the Old Faithful geyser in Yellowstone Park. Figure 7-4 uses six bars that group the data by 10-minute intervals. This histogram shows a general skewed left pattern, but with 222 observations you are cramming an awful lot of data into only six groups; for example, the bar for 75–85 minutes has more than 90 pieces of data in it. You can break it down further than that.

Figure 7-4:
Histogram #1
showing time
between
eruptions for
Old Faithful
geyser
($n = 222$).

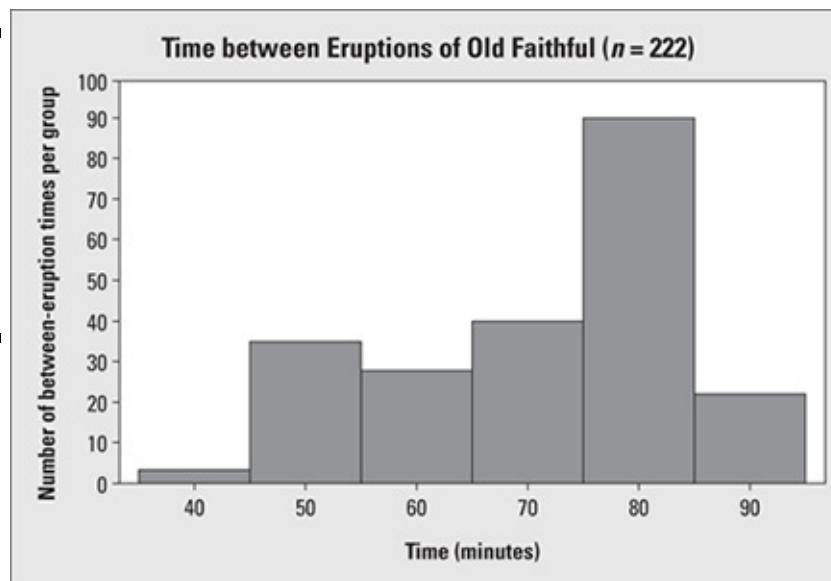


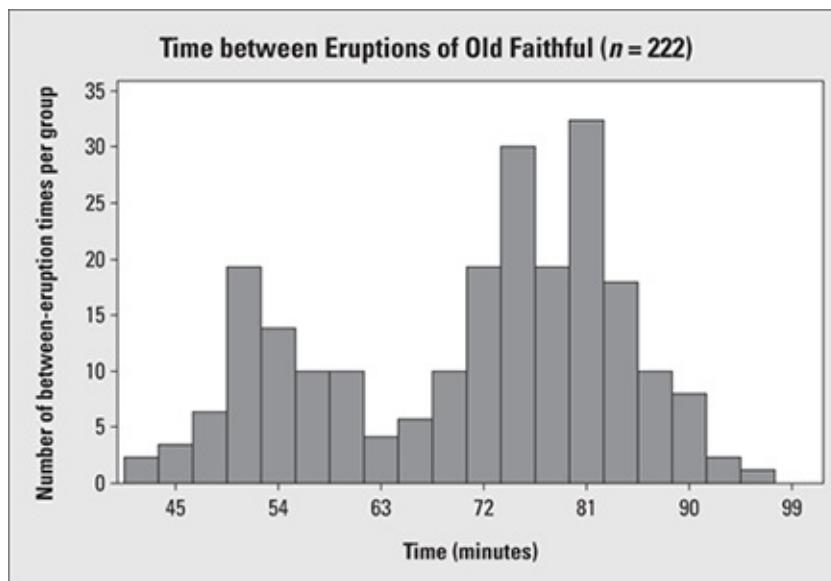
Figure 7-5 is a histogram of the same data set, where the time between eruptions is broken into groups of 3 minutes each, resulting in 19 bars. Notice the distinct pattern in the data that shows up with this histogram which wasn't uncovered in Figure 7-4. You see two distinct peaks in the data; one peak around the 50-minute mark, and one around the 75-minute mark. A data set with two peaks is called *bimodal*; Figure 7-5 shows a clear example.

Looking at Figure 7-5, you can conclude that the geyser has two categories of eruptions; one group that has a shorter waiting time, and another group that has a longer waiting time. Within each group you see the data are fairly close to where the peak is located. Looking at Figure 7-4, you couldn't say that.



If the interval for the groupings of the numerical variable is really small, you see too many bars in the histogram; the data may be hard to interpret because the heights of the bars look more variable than they should be. On the other hand, if the ranges are really large, you see too few bars, and you may miss something interesting in the data.

Figure 7-5:
Histogram #2
showing time
between
eruptions for
Old Faithful
geyser
($n = 222$).



Watching the scale and start/finish lines

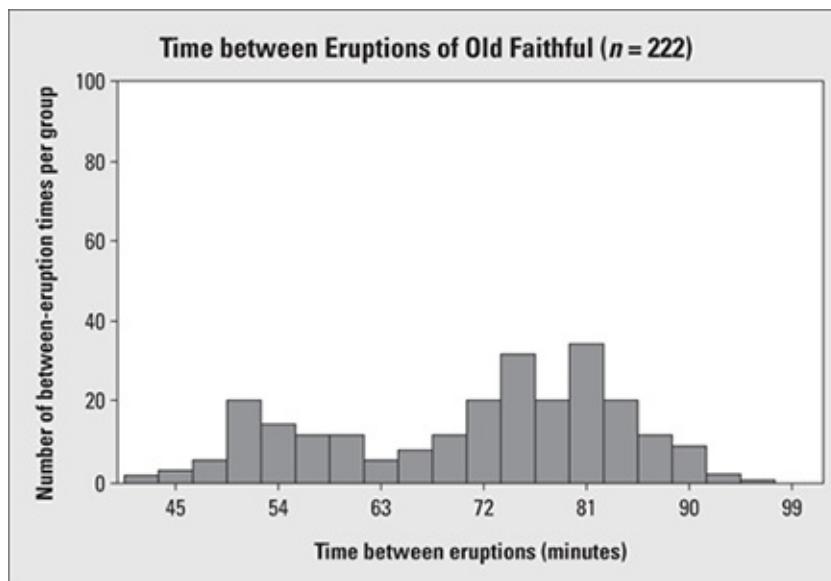
The y-axis of a histogram shows how many individuals are in each group, using counts or percents. A histogram can be misleading if it has a deceptive scale and/or inappropriate starting and ending points on the y-axis.



Watch the scale on the y-axis of a histogram. If it goes by large increments and has an ending point that's much higher than needed, you see a great deal of white space above the histogram. The heights of the bars are squeezed down, making their differences look more uniform than they should. If the scale goes by small increments and ends at the smallest value possible, the bars become stretched vertically, exaggerating the differences in their heights and suggesting a bigger difference than really exists.

An example comparing scales on the vertical (y) axes is shown in Figures 7-5 and 7-6. I took the Old Faithful data (time between eruptions) and made a histogram with vertical increments of 20 minutes, from 0 to 100; see Figure 7-6. Compare this to Figure 7-5, with vertical increments of 5 minutes, from 0 to 35. Figure 7-6 has a lot of white space and gives the appearance that the times are more evenly distributed among the groups than they really are. It also makes the data set look smaller, if you don't pay attention to what's on the y-axis. Of the two graphs, Figure 7-5 is more appropriate.

Figure 7-6:
Histogram #3
of Old Faithful
geyser
eruption
times.



Examining Boxplots

A *boxplot* is a one-dimensional graph of numerical data based on the five-number summary, which includes the minimum value, the 25th percentile (known as Q_1), the median, the 75th percentile (Q_3), and the maximum value. In essence, these five descriptive statistics divide the data set into four parts; each part contains 25% of the data. (See Chapter 5 for a full discussion of the five-number summary.)

Making a boxplot

To make a boxplot, follow these steps:

1. Find the five-number summary of your data set. (Use the steps outlined in Chapter 5.)
2. Create a vertical (or horizontal) number line whose scale includes the numbers in the five-number summary and uses appropriate units of equal distance from each other.
3. Mark the location of each number in the five-number summary just above the number line (for a horizontal boxplot) or just to the right of the number line (for a vertical boxplot).
4. Draw a box around the marks for the 25th percentile and the 75th percentile.
5. Draw a line in the box where the median is located.
6. Determine whether or not outliers are present.

To make this determination, calculate the IQR (by subtracting $Q_3 - Q_1$); then multiply by 1.5. Add this amount to the value of Q_3 and subtract this amount from Q_1 . This gives you a wider boundary around the median than the box does. Any data points that fall outside this boundary are determined to be outliers.

7. If there are no outliers (according to your results of Step 6), draw lines from the upper and lower edges of the box out to the minimum and maximum values in the data set.

8. If there are outliers (according to your results of Step 6), indicate their location on the boxplot with * signs. Instead of drawing a line from the edge of the box all the way to the most extreme outlier, stop the line at the last data value that isn't an outlier.



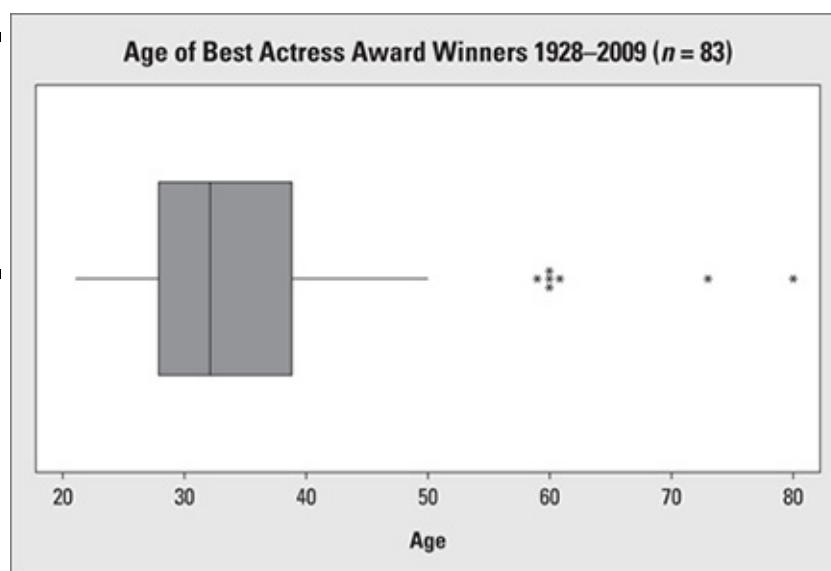
Many if not most software packages indicate outliers in a data set by using an asterisk (*) or star symbol and use the procedure outlined in Step 6 to identify outliers. However, not all packages use these symbols and procedures; check to see what your package does before analyzing your data with a boxplot.

A horizontal boxplot for ages of the Best Actress Oscar award winners from 1928–2009 is shown in Figure 7-7. You can see the numbers separating sections of the boxplot match the five-number summary statistics shown in Figure 7-3.



Boxplots can be vertical (straight up and down) with the values on the axis going from bottom (lowest) to top (highest); or they can be horizontal, with the values on the axis going from left (lowest) to right (highest). The next section shows you how to interpret a boxplot.

Figure 7-7:
Boxplot of
Best Actress
ages (1928–
2009; $n = 83$
actresses).



Interpreting a boxplot

Similar to a histogram (see the section “Interpreting a histogram”), a boxplot can give you information regarding the shape, center, and variability of a data set. Boxplots differ from histograms in terms of their strengths and weaknesses, as you see in the upcoming sections, but one of their biggest strengths is how they handle skewed data.

Checking the shape with caution!

A boxplot can show whether a data set is symmetric (roughly the same on each side when cut down the middle) or skewed (lopsided). A symmetric data set shows the median roughly in the middle of the box. Skewed data show a lopsided boxplot, where the median cuts the box into two unequal pieces. If the longer part of the box is to the right (or above) the median, the data is said to be *skewed right*. If the longer part is to the left (or below) the median, the data is *skewed left*.

As shown in the boxplot of the data in Figure 7-7, the ages are skewed right. The part of the box to the left of the median (representing the younger actresses) is shorter than the part of the box to the right of the median (representing the older actresses). That means the ages of the younger actresses are closer together than the ages of the older actresses. Figure 7-3 shows the descriptive statistics of the data and confirms the right skewness: the median age (33 years) is lower than the mean age (35.69 years).



If one side of the box is longer than the other, it does not mean that side contains more data. In fact, you can't tell the sample size by looking at a boxplot; it's based on percentages, not counts. Each section of the boxplot (the minimum to Q_1 , Q_1 to the median, the median to Q_3 , and Q_3 to the maximum) contains 25% of the data no matter what. If one of the sections is longer than another, it indicates a wider range in the values of data in that section (meaning the data are more spread out). A smaller section of the boxplot indicates the data are more condensed (closer together).



Although a boxplot can tell you whether a data set is symmetric (when the median is in the center of the box), it can't tell you the shape of the symmetry the way a histogram can. For example, Figure 7-8 shows histograms from two different data sets, each one containing 18 values that vary from 1 to 6. The histogram on the left has an equal number of values in each group, and the one on the right has two peaks at 2 and 5. Both histograms show the data are symmetric, but their shapes are clearly different.

Figure 7-8:
Histograms of
two
symmetric
data sets.

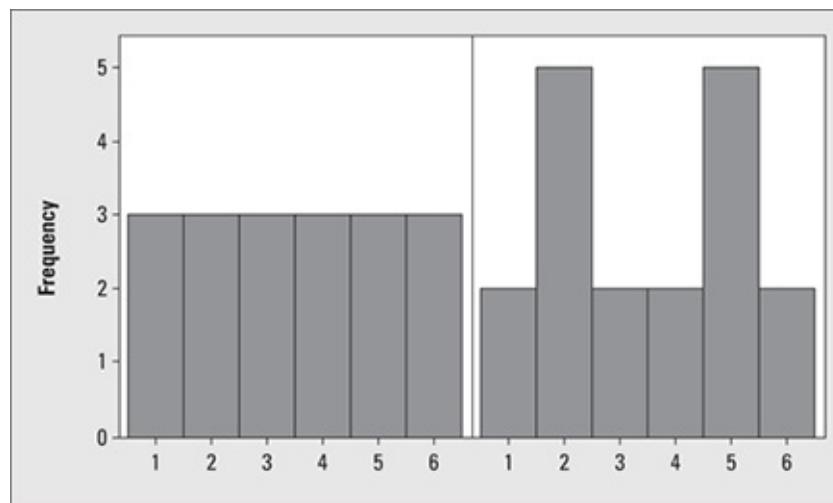
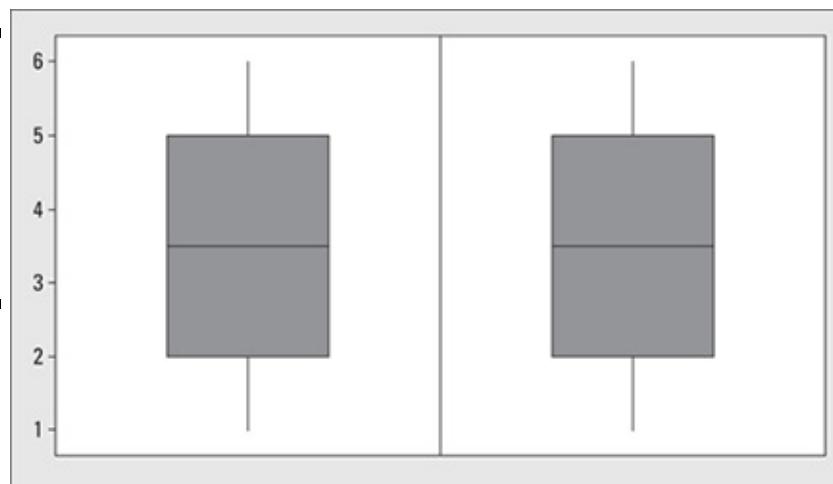


Figure 7-9 shows the corresponding boxplots for these same two data sets; notice they are exactly the same. This is because the data sets both have the same five-number summaries — they’re both symmetric with the same amount of distance between Q_1 , the median, and Q_3 . However, if you just saw the boxplots and not the histograms, you might think the shapes of the two data sets are the same, when indeed they are not.

Figure 7-9:
Boxplots of
the two
symmetric
data sets
from Figure 7-
8.



Despite its weakness in detecting the type of symmetry (you can add in a histogram to your analyses to help fill in that gap), a boxplot has a great upside in that you can identify actual measures of spread and center directly from the boxplot, where on a histogram you can’t. A boxplot is also good for comparing data sets by showing them on the same graph, side by side.



All graphs have strengths and weaknesses; it’s always a good idea to show more than one graph of your data for that reason.

Measuring variability with IQR

Variability in a data set that is described by the five-number summary is measured by the interquartile range (*IQR*). The *IQR* is equal to $Q_3 - Q_1$, the difference between the 75th

percentile and the 25th percentile (the distance covering the middle 50% of the data). The larger the *IQR*, the more variable the data set is.

From Figure 7-3, the variability in age of the Best Actress winners as measured by the *IQR* is $Q_3 - Q_1 = 39 - 28 = 11$ years. Of the group of actresses whose ages were closest to the median, half of them were within 11 years of each other when they won their awards.



Notice that the *IQR* ignores data below the 25th percentile or above the 75th, which may contain outliers that could inflate the measure of variability of the entire data set. So if data is skewed, the *IQR* is a more appropriate measure of variability than the standard deviation.

Picking out the center using the median

The median, part of the five-number summary, is shown by the line that cuts through the box in the boxplot. This makes it very easy to identify. The mean, however, is not part of the boxplot and can't be determined accurately by just looking at the boxplot.

You don't see the mean on a boxplot because boxplots are based completely on percentiles. If data are skewed, the median is the most appropriate measure of center. Of course you can calculate the mean separately and add it to your results; it's never a bad idea to show both.

Investigating Old Faithful's boxplot

The relevant descriptive statistics for the Old Faithful geyser data are found in Figure 7-10.

Figure 7-10:
Descriptive statistics for Old Faithful data.

Descriptive Statistics: Time between Eruptions										
Variable	Total	Count	Mean	StDev	Minimum	Q1	Median	Q3	Maximum	IQR
Time between	222	71.009	12.799	42.000	60.000	75.000	81.000	95.000	21.000	

You can predict from the data set that the shape will be skewed left a bit because the mean is lower than the median by about 4 minutes. The *IQR* is $Q_3 - Q_1 = 81 - 60 = 21$ minutes, which shows the amount of overall variability in the time between eruptions; 50% of the eruptions are within 21 minutes of each other.

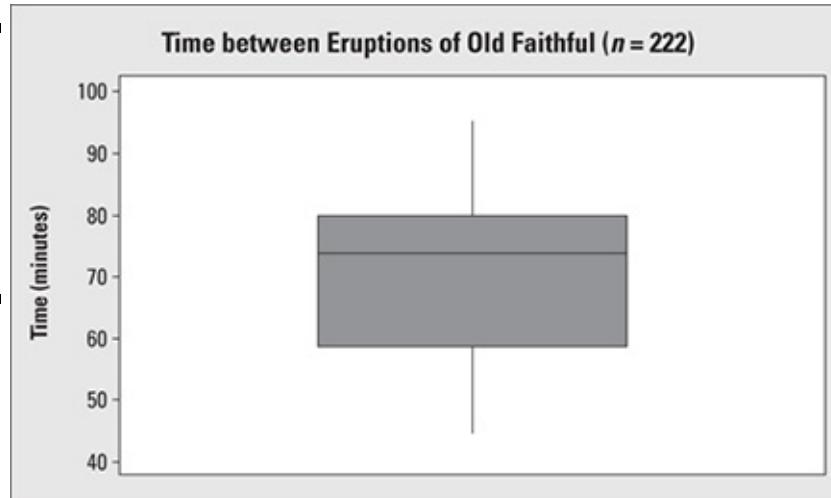
A vertical boxplot for length of time between eruptions of the Old Faithful geyser is shown in Figure 7-11. You confirm that the data are skewed left because the lower part of

the box (where the small values are) is longer than the upper part of the box.

You see the values of the boxplot in Figure 7-11 that mark the five-number summary and the information shown in Figure 7-10, including the *IQR* of 21 minutes to measure variability. The center as marked by the median is 75 minutes; this is a better measure of center than the mean (71 minutes), which is driven down a bit by the left skewed values (the few that are shorter times than the rest of the data).

Looking at the boxplot (Figure 7-11), you see there are no outliers denoted by stars. However, note that the boxplot doesn't pick up on the bimodal shape of the data that you see in Figure 7-5. You need a good histogram for that.

Figure 7-11:
Boxplot of
eruption
times for Old
Faithful
geyser
($n = 222$).



Denoting outliers

Looking at the boxplot in Figure 7-7 for the Best Actress ages data, you see a set of outliers (seven in all) on the right side of the data set, marked by a group of stars (as described in Step 8 in the earlier section “Making a boxplot”). Three of the stars lie on top of one another because three actresses were the same age, 61, when they won their Oscars.

You verify these outliers by applying the rule described in Step 6 of the section “Making a Boxplot.” The *IQR* is 11 (from Figure 7-3), so you take $11 * 1.5 = 16.5$ years. Add this amount to Q_3 and you get $39 + 16.5 = 55.5$ years; subtracting this amount from Q_1 you get $28 - 16.5 = 11.5$ years. So an actress whose age was below 11.5 years (that is, 11 years old and under) or above 55.5 years (that is, 56 years old or over) is considered to be an outlier.

Of course, the lower end of this boundary (11.5 years) isn't relevant because the youngest actress was 21 (Figure 7-3 shows the minimum is 21). So you know there aren't any outliers on the low end of this data set.

However, seven outliers are on the high end of the data set, where the 56-and-over actresses' ages are. Table 7-2 shows the information on all seven outliers in the Best

Table 7-2**Best Actress Winners with Ages Designated as Outliers**

Year	Name	Age	Movie
1967	Katharine Hepburn	60	<i>Guess Who's Coming to Dinner</i>
1968	Katharine Hepburn	61	<i>The Lion in Winter</i>
1985	Geraldine Page	61	<i>Trip to Bountiful</i>
2006	Helen Mirren	61	<i>The Queen</i>
1931	Marie Dressler	62	<i>Min and Bill</i>
1981	Katharine Hepburn	74	<i>On Golden Pond</i>
1989	Jessica Tandy	81	<i>Driving Miss Daisy</i>

The youngest of the outliers is 60 years old (Katharine Hepburn, 1967). Just to compare, the next youngest age in the data set is 49 (Susan Sarandon, 1995). This indicates a clear break in this data set.

Making mistakes when interpreting a boxplot

It's a common mistake to associate the size of the box in a boxplot with the amount of data in the data set. Remember that each of the four sections shown in the boxplot contains an equal percentage (25%) of the data; the boxplot just marks off the places in the data set that separate those sections.



In particular, if the median splits the box into two unequal parts, the larger part contains data that's more variable than the other part, in terms of its range of values. However, there is still the same amount of data (25%) in the larger part of the box as there is in the smaller part.

Another common error involves sample size. A boxplot is a one-dimensional graph with only one axis representing the variable being measured. There is no second axis that tells you how many data points are in each group. So if you see two boxplots side-by-side and one of them has a very long box and the other has a very short one, don't conclude that the longer one has more data in it. The length of the box represents the variability in the data, not the number of data values.



When viewing or making a boxplot, always make sure the sample size (n) is included as part of the title. You can't figure out the sample size otherwise.

Tackling Time Charts

A *time chart* (also called a *line graph*) is a data display used to examine trends in data over time (also known as time series data). Time charts show time on the *x*-axis (for example, by month, year, or day) and the values of the variable being measured on the *y*-axis (like birth rates, total sales, or population size). Each point on the time chart summarizes all the data collected at that particular time; for example, the average of all pepper prices for January or the total revenue for 2010.

Interpreting time charts



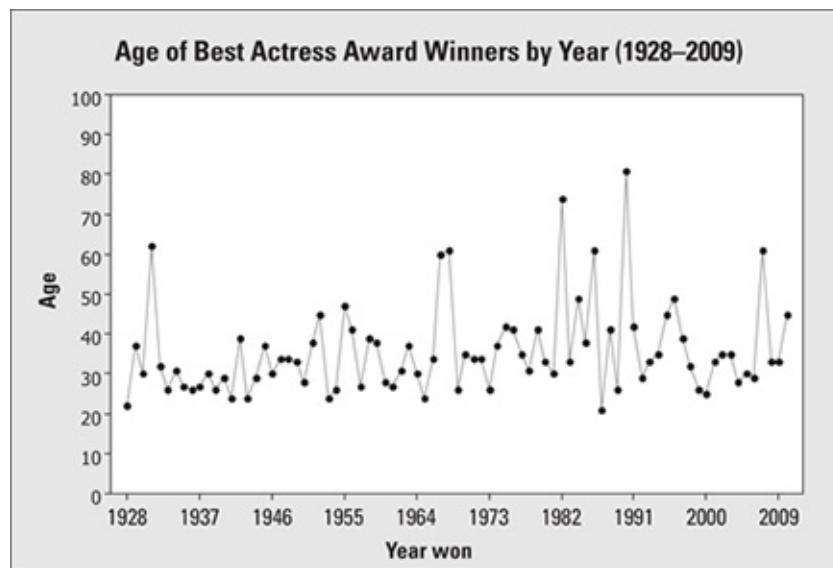
To interpret a time chart, look for patterns and trends as you move across the chart from left to right.

The time chart in Figure 7-12 shows the ages of the Best Actress winners, in order of year won, from 1928–2009. Each dot indicates the age of a single actress, the one that won the Oscar that year. You see a bit of a cyclical pattern across time; that is, the ages go up, down, up, down, up, down with at least some regularity. It's hard to say what may be going on here; many variables go into determining an Oscar winner, including the type of movie, type of female role, mood of the voters, and so forth, and some of these variables may have a cyclical pattern to them.

Figure 7-12 also shows a very faint trend in age that is tending uphill; indicating that the Best Actress Award winners may be winning their awards increasingly later in life. Again, I wouldn't make too many assumptions from this result because the data has a great deal of variability.

As far as variability goes, you see that the ages represented by the dots do fluctuate quite a bit on the *y*-axis (representing age); all the dots basically fall between 20 and 80 years, with most of them between 25 and 45 years, I'd say. This goes along with the descriptive statistics found in Figure 7-3.

Figure 7-12:
Time Chart #1
for ages of
Best Actress
Academy
Award
winners,
1928–2009.



Understanding variability: Time charts versus histograms



Variability in a histogram should not be confused with variability in a time chart. If values change over time, they're shown on a time chart as highs and lows, and many changes from high to low (over time) indicate lots of variability. So a flat line on a time chart indicates no change and no variability in the values across time. For example, if the price of a product stays the same for 12 months in a row, the time chart for price would be flat.

But when the heights of a histogram's bars appear flat, the data is spread out uniformly across all the groups, indicating a great deal of variability in the data. (For an example, refer to Figure 7-2a.)

Spotting misleading time charts

As with any graph, you have to evaluate the units of the numbers being plotted. For example, it's misleading to chart the *number* of crimes over time, rather than the crime *rate* (crimes per capita) — because the population size of a city changes over time, crime rate is the appropriate measure. Make sure you understand what numbers are being graphed and examine them for fairness and appropriateness.

Watching the scale and start/end points

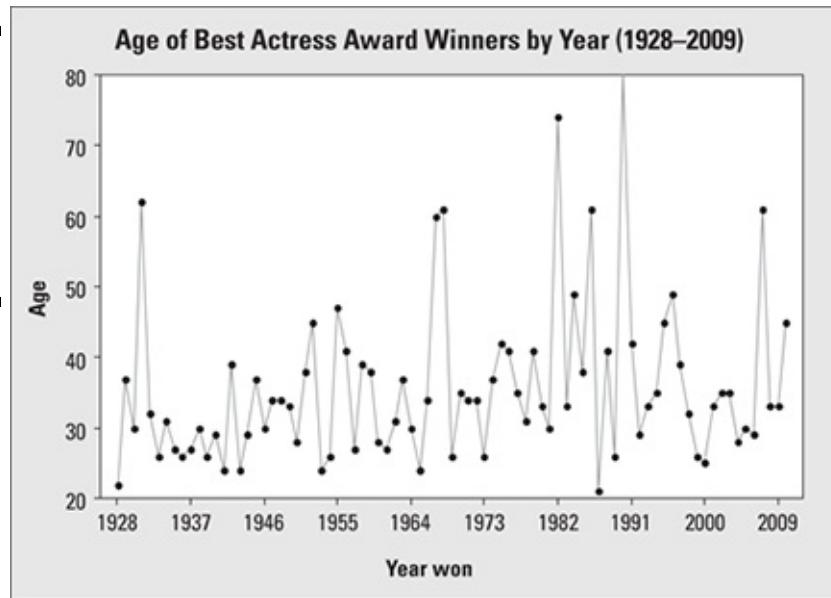
The scale on the vertical axis can make a big difference in the way the time chart looks. Refer to Figure 7-12 to see my original time chart of the ages for the Best Actress Academy Award winners from 1928–2009 in increments of 10 years. You see a fair amount of

variability, as discussed previously.

In Figure 7-12, the starting and ending points on the vertical axis are 0 to 100, which creates a little bit of extra white space on the top and bottom of the picture. I could have used 10 and 90 as my start/end points, but this graph looks reasonable.

Now what happens if I change the vertical axis? Figure 7-13 shows the same data, with start/end points of 20 and 80. The increments of 10 years appear longer than the increments of 10 years shown in Figure 7-12. Both of these changes in the graph exaggerate the differences in ages even more.

Figure 7-13:
Time Chart #2
for ages of
Best Actress
Oscar Award
winners,
1928–2009.



How do you decide which graph is the best one for your data? There is no perfect graph; there is no right or wrong answer; but there are limits. You can quickly spot problems just by zooming in on the scale and start/end points.

Simplifying excess data

A time chart of the time between eruptions for the Old Faithful data is shown in Figure 7-14. You see 222 dots on this graph; each one represents the time between one eruption and the next, for every eruption during a 16-day period.

This figure looks very complex; data are everywhere, there are too many points to really see anything, and you can't find the forest for the trees. There is such a thing as having too much data, especially nowadays when you can measure data continuously and meticulously using all kinds of advanced technology. I'm betting they didn't have a student standing by the geyser recording eruption times on a clipboard, for example!

To get a clearer picture of the Old Faithful data, I combined all the observations from a single day and found its mean; I did this for all 16 days, and then I plotted all the means on a time chart in order. This reduced the data from 222 points to 16 points. The time

chart is shown in Figure 7-15.

From this time chart I see a little bit of a cyclical pattern to the data; every day or two it appears to shift from short times between eruptions to longer times between eruptions. While these changes are not definitive, it does provide important information for scientists to follow up on when studying the behavior of geysers like Old Faithful.

Figure 7-14:
Time chart showing time between eruptions for Old Faithful Geyser ($n = 222$ consecutive observations).

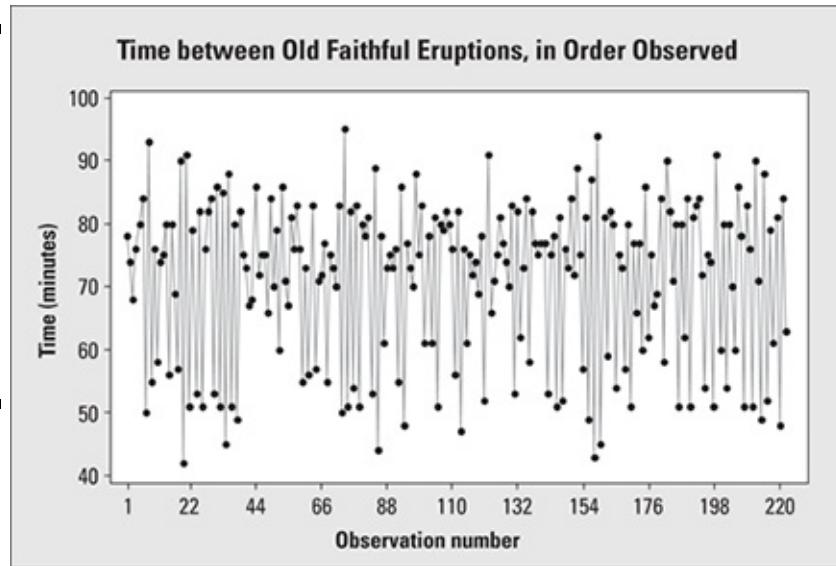
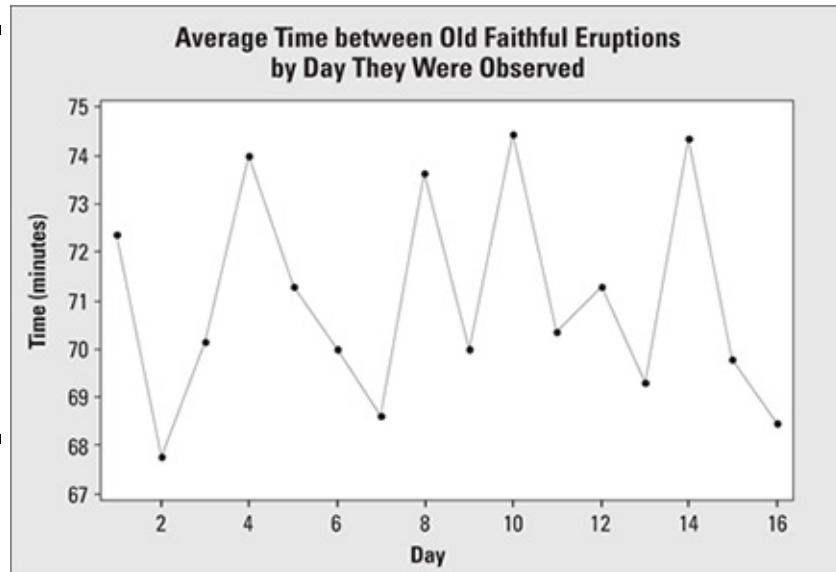


Figure 7-15:
Time chart showing daily average time between eruptions for Old Faithful geyser ($n = 16$ consecutive days).



A time chart condenses all the data for one unit of time into a single point. By contrast, a histogram displays the entire sample of data that was collected at that one unit of time. For example, Figure 7-15 shows the daily average time between eruptions for 16 days. For any given day, you can make a histogram of all the eruptions observed on that particular day. Displaying a time chart of average times over 16 days accompanied by a histogram summarizing all the eruptions for a particular day would be a great one-two punch.

Evaluating time charts

Here is a checklist for evaluating time charts, with a couple more thoughts added in:

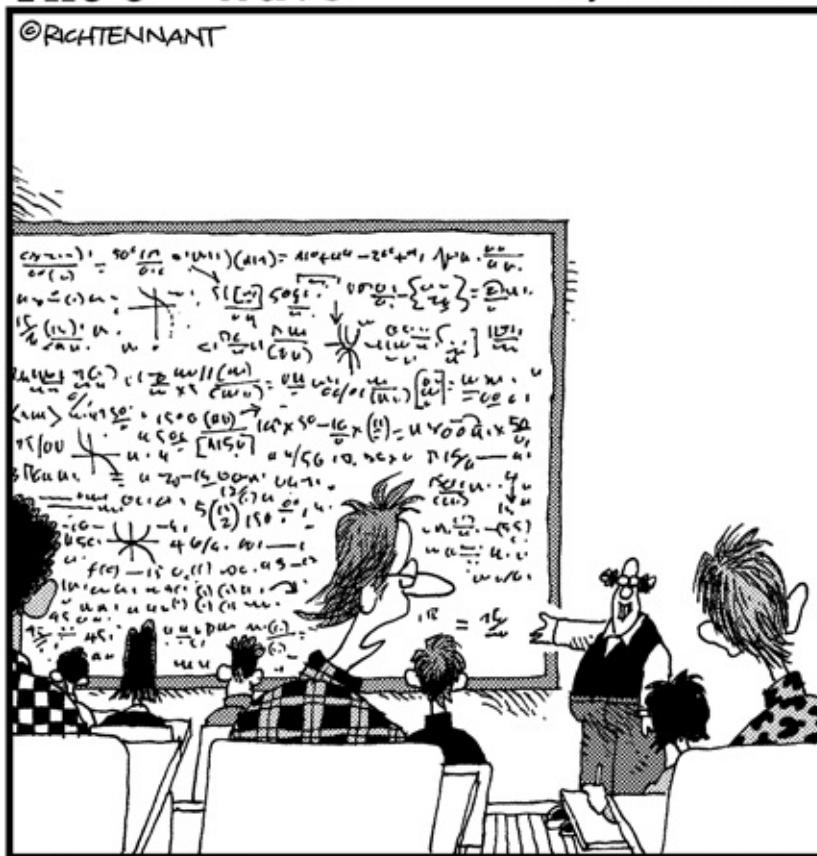
- ✓ Examine the scale and start/end points on the vertical axis (the one showing the values of the data). Large increments and/or lots of white space make differences look less dramatic; small increments and/or a plot that totally fills the page exaggerate differences.
- ✓ If the amount of data you have is overwhelming, consider boiling it down by finding means/medians for blocks of time and plotting those instead.
- ✓ Watch for gaps in the timeline on a time chart. For example, it's misleading to show equally spaced points on the horizontal (time) axis for 1990, 2000, 2005, and 2010. This happens when years are just treated like labels, rather than real numbers.
- ✓ As with any graph, take the units into account; be sure they're appropriate for comparison over time. For example, are dollar amounts adjusted for inflation? Are you looking at number of crimes, or the crime rate?

Part III

Distributions and the Central Limit Theorem

The 5th Wave

By Rich Tennant



"This guy writes an equation for over 20 minutes,
and he has the nerve to say, 'Voilà'?"

In this part . . .

Statisticians study populations; that's their bread and butter. They measure, count, or classify characteristics of a population (using random variables); find probabilities and proportions; and create (or estimate) numerical summaries for the population (that is, parameters for the population). Sometimes you know a great deal about a population from the start; sometimes it's hazier. This part studies populations under both scenarios.

If a population fits a specific distribution, tools are available for studying it. In Chapters 8 through 10, you see three commonly used distributions: the binomial distribution (for categorical data) and the normal and t-distributions (for numerical data).

If the specifics about a population are unknown (as happens most of the time), you take a sample and generalize its results to the population. However, sample results vary, and you need to take that into account. In Chapter 11 you investigate sample variability,

measure the precision of your sample results, and find probabilities for their likelihood. From there you'll be able to properly estimate parameters and test claims made about them, but that's another Part — IV, to be exact.

Chapter 8

Random Variables and the Binomial Distribution

In This Chapter

- ▶ Identifying a binomial random variable
 - ▶ Finding probabilities using a formula or table
 - ▶ Calculating the mean and variance
-

Scientists and engineers often build models for the phenomena they are studying to make predictions and decisions. For example, where and when is this hurricane going to hit when it makes landfall? How many accidents will occur at this intersection this year if it's not redone? Or, what will the deer population be like in a certain region five years from now?

To answer these questions, scientists (usually working with statisticians) define a characteristic they are measuring or counting (such as number of intersections, location and time when a hurricane hits, population size, and so on) and treat it as a variable that changes in some random way, according to a certain pattern. They cleverly call them — you guessed it — random variables. In this chapter, you find out more about random variables, their types and characteristics, and why they are important. And you look at the details of one of the most common random variables: the binomial.

Defining a Random Variable

A *random variable* is a characteristic, measurement, or count that changes randomly according to a certain set or pattern. Its notation is X , Y , Z , and so on. In this section, you see how different random variables are characterized and how they behave in the long term in terms of their means and standard deviations.



In math you have variables like X and Y that take on certain values depending on the problem (for example, the width of a rectangle), but in statistics the variables change in a random way. By *random*, statisticians mean that you don't know exactly what the next outcome will be but you do know that certain outcomes happen more frequently than others; everything's not 50-50. (Like when I try to shoot baskets; it's definitely not a 50% chance I'll make one and 50% chance I'll miss. It's more like 5% chance of making it and a 95% chance of missing it.) You can use that information to better study data and populations and make good decisions. (For example, don't put

me in your basketball game to shoot free throws.)

Data have different types: categorical and numerical (see Chapter 4). While both types of data are associated with random variables, I discuss only numerical random variables here (this falls in line with most intro stat courses as well). For information on analyzing categorical variables, see Chapters 6 and 19.

Discrete versus continuous

Numerical random variables represent counts and measurements. They come in two different flavors: discrete and continuous, depending on the type of outcomes that are possible.

- ✓ **Discrete random variables:** If the possible outcomes of a random variable can be listed out using whole numbers (for example, 0, 1, 2 . . . , 10; or 0, 1, 2, 3), the random variable is *discrete*.
- ✓ **Continuous random variables:** If the possible outcomes of a random variable can only be described using an interval of real numbers (for example, all real numbers from zero to infinity), the random variable is *continuous*.

Discrete random variables typically represent counts — for example, the number of people who voted yes for a smoking ban out of a random sample of 100 people (possible values are 0, 1, 2, . . . , 100); or the number of accidents at a certain intersection over one year's time (possible values are 0, 1, 2, . . .).



Discrete random variables have two classes: finite and countably infinite. A discrete random variable is *finite* if its list of possible values has a fixed (finite) number of elements in it (for example, the number of smoking ban supporters in a random sample of 100 voters has to be between 0 and 100). One very common finite random variable is the binomial, which is discussed in this chapter in detail.

A discrete random variable is *countably infinite* if its possible values can be specifically listed out but they have no specific end. For example, the number of accidents occurring at a certain intersection over a 10-year period can take on possible values: 0, 1, 2, . . . (you know they end somewhere but you can't say where, so you list them all).

Continuous random variables typically represent measurements, such as time to complete a task (for example 1 minute 10 seconds, 1 minute 20 seconds, and so on) or the weight of a newborn. What separates continuous random variables from discrete ones is that they are *uncountably infinite*; they have too many possible values to list out or to count and/or they can be measured to a high level of precision (such as the level of smog in the air in Los Angeles on a given day, measured in parts per million).

Examples of commonly used continuous random variables can be found in Chapter 9 (the normal distribution) and Chapter 10 (the *t*-distribution).

Probability distributions

A discrete random variable X can take on a certain set of possible outcomes, and each of those outcomes has a certain probability of occurring. The notation used for any specific outcome is a lowercase x . For example, say you roll a die and look at the outcome. The random variable X is the outcome of the die (which takes on possible values of 1, 2, . . . , 6). Now if you roll the die and get a 1, that's a specific outcome, so you write " $x = 1$."

The probability of any specific outcome occurring is denoted $p(x)$, which you pronounce " p of x ." It signifies the probability that the random variable X takes on a specific value, which you call "little x ." For example, to denote the probability of getting a 1 on a die, you write $p(1)$.



Statisticians use an uppercase X when they talk about random variables in their general form; for example, "Let X be the outcome of the roll of a single die." They use lowercase x when they talk about specific outcomes of the random variable, like $x = 1$ or $x = 2$.

A list or function showing all possible values of a discrete random variable, along with their probabilities, is called a *probability distribution*, $p(x)$. For example, when you roll a single die, the possible outcomes are 1, 2, 3, 4, 5, and 6, and each has a probability of $1/6$ (if the die is fair). As another example, suppose 40% of renters living in an apartment complex own one dog, 7% own two dogs, 3% own three dogs, and 50% own zero dogs. For X = the number of dogs owned, the probability distribution for X is shown in Table 8-1.

Table 8-1 Probability Distribution for X = Number of Dogs Owned by Apartment Renters

x	$p(x)$
0	0.50
1	0.40
2	0.07
3	0.03

The mean and variance of a discrete random variable

The *mean* of a random variable is the average of all the outcomes you would expect in the long term (over all possible samples). For example, if you roll a die a billion times and record the outcomes, the average of those outcomes is 3.5. (Each outcome happens

with equal chance, so you average the numbers 1 through 6 to get 3.5.) However, if the die is loaded and you roll a 1 more often than anything else, the average outcome from a billion rolls is closer to 1 than to 3.5.



The notation for the mean of a random variable X is μ_x or μ (pronounced “mu sub x ”; or just “mu x ”). Because you are looking at all the outcomes in the long term, it’s the same as looking at the mean of an entire population of values, which is why you denote it μ_x and not \bar{x} . (The latter represents the mean of a *sample* of values [see Chapter 5].) You put the X in the subscript to remind you that the variable this mean belongs to is the X variable (as opposed to a Y variable or some other letter).

The *variance* of a random variable is roughly interpreted as the average squared distance from the mean for all the outcomes you would get in the long term, over all possible samples. This is the same as the variance of the population of all possible values. The notation for variance of a random variable X is σ_x^2 or σ^2 . You say “sigma sub x , squared” or just “sigma squared.”

The standard deviation of a random variable X is the square root of the variance, denoted by σ_x or σ (say “sigma x ” or just “sigma”). It roughly represents the average distance from the mean.

Just like for the mean, you use the Greek notation to denote the variance and standard deviation of a random variable. The English notation s^2 and s represent the variance and standard deviation of a *sample* of individuals, not the entire population (see Chapter 5).



The variance is in square units, so it can’t be easily interpreted. You use standard deviation for interpretation because it is in the original units of X . The standard deviation can be roughly interpreted as the average distance away from the mean.

Identifying a Binomial

The most well-known and loved discrete random variable is the binomial. *Binomial* means *two names* and is associated with situations involving two outcomes; for example yes/no, or success/failure (hitting a red light or not, developing a side effect or not). This section focuses on the binomial random variable — when you can use it, finding probabilities for it, and finding its mean and variance.

A random variable is binomial (that is, it has a binomial distribution) if the following four conditions are met:

1. There are a fixed number of trials (n).

2. Each trial has two possible outcomes: success or failure.
3. The probability of success (call it p) is the same for each trial.
4. The trials are independent, meaning the outcome of one trial doesn't influence that of any other.

Let X equal the total number of successes in n trials; if all four conditions are met, X has a binomial distribution with probability of success (on each trial) equal to p .

The lowercase p here stands for the probability of getting a success on one single (individual) trial. It's not the same as $p(x)$, which means the probability of getting x successes in n trials.

Checking binomial conditions step by step

You flip a fair coin 10 times and count the number of heads (X). Does X have a binomial distribution? You can check by reviewing your responses to the questions and statements in the list that follows:

1. Are there a fixed number of trials?

You're flipping the coin 10 times, which is a fixed number. Condition 1 is met, and $n = 10$.

2. Does each trial have only two possible outcomes — success or failure?

The outcome of each flip is either heads or tails, and you're interested in counting the number of heads. That means success = heads, and failure = tails. Condition 2 is met.

3. Is the probability of success the same for each trial?

Because the coin is fair, the probability of success (getting a head) is $p = 1/2$ for each trial. You also know that $1 - 1/2 = 1/2$ is the probability of failure (getting a tail) on each trial. Condition 3 is met.

4. Are the trials independent?

You assume the coin is being flipped the same way each time, which means the outcome of one flip doesn't affect the outcome of subsequent flips. Condition 4 is met.

Because the random variable X (the number of successes [heads] that occur in 10 trials [flips]) meets all four conditions, you conclude it has a binomial distribution with $n = 10$ and $p = 1/2$.

But not every situation that appears binomial actually is. Read on to see some examples of what I mean.

No fixed number of trials

Suppose that you're going to flip a fair coin until you get four heads and you'll count how many flips it takes to get there; in this case X = number of flips. This certainly sounds like a binomial situation: Condition 2 is met because you have success (heads) and failure (tails) on each flip; condition 3 is met with the probability of success (heads) being the same (0.5) on each flip; and the flips are independent, so condition 4 is met.

However, notice that X isn't counting the number of heads, it counts the number of trials needed to get 4 heads. The number of successes (X) is fixed rather than the number of trials (n). Condition 1 is not met, so X does not have a binomial distribution in this case.

More than success or failure

Some situations involve more than two possible outcomes, yet they can appear to be binomial. For example, suppose you roll a fair die 10 times and let X be the outcome of each roll (1, 2, 3, . . . , 6). You have a series of $n = 10$ trials, they are independent, and the probability of each outcome is the same for each roll. However, on each roll you're recording the outcome on a six-sided die, a number from 1 to 6. This is not a success/failure situation, so condition 2 is not met.

However, depending on what you're recording, situations originally having more than two outcomes can fall under the binomial category. For example, if you roll a fair die 10 times and each time you record whether or not you get a 1, then condition 2 is met because your two outcomes of interest are getting a 1 ("success") and not getting a 1 ("failure"). In this case, p (the probability of success) = 1/6, and 5/6 is the probability of failure. So if X is counting the number of 1s you get in 10 rolls, X is a binomial random variable.

Trials are not independent

The independence condition is violated when the outcome of one trial affects another trial. Suppose you want to know opinions of adults in your city regarding a proposed casino. Instead of taking a random sample of, say, 100 people, to save time you select 50 married couples and ask each of them what their opinion is. In this case it's reasonable to say couples have a higher chance of agreeing on their opinions than individuals selected at random, so the independence condition 4 is not met.

Probability of success (p) changes

You have 10 people — 6 women and 4 men — and you want to form a committee of 2 people at random. Let X be the number of women on the committee of 2. The chance of selecting a woman at random on the first try is 6/10. Because you can't select this same woman again, the chance of selecting another woman is now 5/9. The value of p has changed, and condition 3 is not met.



If the population is very large (for example all U.S. adults), p still changes every time you choose someone, but the change is negligible, so you don't worry about it. You still say the trials are independent with the same probability of success, p . (Life is so much easier that way!)

Finding Binomial Probabilities Using a Formula

After you identify that X has a binomial distribution (the four conditions from the section “Checking binomial conditions step by step” are met), you’ll likely want to find probabilities for X . The good news is that you don’t have to find them from scratch; you get to use established formulas for finding binomial probabilities, using the values of n and p unique to each problem. Probabilities for a binomial random variable X can be found using the following formula for $p(x)$:

$$\binom{n}{x} p^x (1-p)^{n-x}$$

where

- ✓ n is the fixed number of trials.
- ✓ x is the specified number of successes.
- ✓ $n - x$ is the number of failures.
- ✓ p is the probability of success on any given trial.
- ✓ $1 - p$ is the probability of failure on any given trial. (**Note:** Some textbooks use the letter q to denote the probability of failure rather than $1 - p$.)

These probabilities hold for any value of X between 0 (lowest number of possible successes in n trials) and n (highest number of possible successes).



The number of ways to rearrange x successes among n trials is called “ n choose x ,” and the notation is $\binom{n}{x}$. It’s important to note that this math expression is not a fraction; it’s math shorthand to represent the number of ways to do these types of rearrangements.

In general, to calculate “ n choose x ,” you use the following formula:

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

The notation $n!$ stands for *n-factorial*, the number of ways to rearrange n items. To calculate $n!$, you multiply $n(n - 1)(n - 2) \dots (2)(1)$. For example $5!$ is $5(4)(3)(2)(1) = 120$; $2!$ is $2(1) = 2$; and $1!$ is 1. By convention, $0!$ equals 1.

Suppose you have to cross three traffic lights on your way to work. Let X be the number of red lights you hit out of the three. How many ways can you hit two red lights on your way to work? Well, you could hit a green one first, then the other two red; or you could hit the green one in the middle and have red ones for the first and third lights, or you could hit red first, then another red, then green. Letting G = green and R=red, you can write these three possibilities as: GRR, RGR, RRG. So you can hit two red lights on your way to work in three ways, right?

Check the math. In this example, a “trial” is a traffic light; and a “success” is a red light. (I know, that seems weird, but a success is whatever you are interested in counting, good or bad.) So you have $n = 3$ total traffic lights, and you’re interested in the situation where

you get $x = 2$ red ones. Using the fancy notation, $\binom{3}{2}$ means “3 choose 2” and stands for the number of ways to rearrange 2 successes in 3 trials.

To calculate “3 choose 2,” you do the following:

$$\binom{3}{2} = \frac{3!}{2!(3-2)!} = \frac{3(2)(1)}{[(2)(1)](1)} = \frac{6}{2} = 3$$

This confirms the three possibilities listed for getting two red lights.

Now suppose the lights operate independently of each other and each one has a 30% chance of being red. Suppose you want to find the probability distribution for X . (That is, a list of all possible values of X — 0,1,2,3 — and their probabilities.)

Before you dive into the calculations, you first check the four conditions (from the section “Checking binomial conditions step by step”) to see if you have a binomial situation here. You have $n = 3$ trials (traffic lights) — check. Each trial is success (red light) or failure (yellow or green light; in other words, “non-red” light) — check. The lights operate independently, so you have the independent trials taken care of, and because each light is red 30% of the time, you know $p = 0.30$ for each light. So X = number of red traffic lights has a binomial distribution. To fill in the nitty gritties for the formulas, $1 - p$ = probability of a non-red light = $1 - 0.30 = 0.70$; and the number of non-red lights is $3 - X$.

Using the formula for $p(x)$, you obtain the probabilities for $x = 0, 1, 2$, and 3 red lights:

$$p(0) = \binom{3}{0} 0.30^0 (1-0.30)^{3-0} =$$

$$\frac{3!}{0!(3-0)!} (0.30)^0 (0.70)^3 = 1(1)(0.343) = 0.343$$

$$p(1) = \binom{3}{1} 0.30^1 (1-0.30)^{3-1} =$$

$$\frac{3!}{1!(3-1)!} (0.30)^1 (0.70)^2 = 3(0.30)(0.49) = 0.441$$

$$p(2) = \binom{3}{2} 0.30^2 (1-0.30)^{3-2} =$$

$$\frac{3!}{2!(3-2)!} (0.30)^2 (0.70)^1 = 3(0.09)(0.70) = 0.189$$

; and

$$p(3) = \binom{3}{3} 0.30^3 (1-0.30)^{3-3} =$$

$$\frac{3!}{3!(3-3)!} (0.30)^3 (0.70)^0 = 1(0.027)(1) = 0.027$$

The final probability distribution for X is shown in Table 8-2. Notice these probabilities all sum to 1 because every possible value of X is listed and accounted for.

Table 8-2 **Probability Distribution for X = Number of Red Traffic Lights ($n = 3$, $p = 0.30$)**

X	$p(x)$
0	0.343
1	0.441
2	0.189
3	0.027

Finding Probabilities Using the Binomial Table

The previous section deals with values of n that are pretty small, but you may wonder how you are going to handle the formula for calculating binomial probabilities when n gets large. No worries! A large range of binomial probabilities are provided in the binomial table in the appendix. Here's how to use it:

Within the binomial table you see several mini-tables; each one corresponds with a different n for a binomial ($n = 1, 2, 3, \dots, 15$, and 20 are available). Each mini-table has rows and columns. Running down the side of any mini-table, you see all the possible values of X from 0 through n , each with its own row. The columns of the binomial table represent various values of p from 0.10 through 0.90.

Finding probabilities for specific values of X

To use the binomial table in the appendix to find probabilities for X = total number of successes in n trials where p is the probability of success on any individual trial, follow these steps:

- 1. Find the mini-table associated with your particular value of n (the number of trials).**
- 2. Find the column that represents your particular value of p (or the one closest to it, if appropriate).**
- 3. Find the row that represents the number of successes (x) you are interested in.**
- 4. Intersect the row and column from Steps 2 and 3.** This gives you the probability for x successes, written as $p(x)$.

For the traffic light example from “Finding Binomial Probabilities Using a Formula,” you can use the binomial table (Table A-3 in the appendix) to verify the results found by the binomial formula shown back in Table 8-2. Go to the mini-table where $n = 3$ and look in the column where $p = 0.30$. You see four probabilities listed for this mini-table: 0.343, 0.441, 0.189, and 0.027; these are the probabilities for $X = 0, 1, 2$, and 3 red lights, respectively, matching those from Table 8-2.

Finding probabilities for X greater-than, less-than, or between two values

The binomial table (Table A-3 in the appendix) shows probabilities for X being equal to any value from 0 to n , for a variety of p s. To find probabilities for X being less-than, greater-than, or between two values, just find the corresponding values in the table and add their probabilities. For the traffic light example, you count the number of times (X) that you hit a red light (out of 3 possible lights). Each light has a 0.30 chance of being red, so you have a binomial distribution with $n = 3$ and $p = 0.30$. If you want the probability that you hit more than one red light, you find $p(x > 1)$ by adding $p(2) + p(3)$ from Table A-3 to get $0.189 + 0.027 = 0.216$.

The probability that you hit between 1 and 3 (inclusive) red lights is $p(1 \leq x \leq 3) = 0.441 + 0.189 + 0.027 = 0.657$.



You have to distinguish between a *greater-than* ($>$) and a *greater-than-or-equal-to* (\geq) probability when working with discrete random variables. Repackaging the previous two examples, you see $p(x > 1) = 0.216$ but $p(x \geq 1) = 0.657$. This is a non-issue for continuous random variables (see Chapter 9).



Other phrases to remember: *at least* means that number or higher, and *at most* means that number or lower. For example, the probability that X is at least 2 is $p(x \geq 2)$; the probability that X is at most 2 is $p(x \leq 2)$.

Checking Out the Mean and Standard Deviation of the Binomial

Because the binomial distribution is so commonly used, statisticians went ahead and did all the grunt work to figure out nice, easy formulas for finding its mean, variance, and standard deviation. (That is, they've already applied the methods from the section "Defining a Random Variable" to the binomial distribution formulas, crunched everything out, and presented the results to us on a silver platter — don't you love it when that happens?) The following results are what came out of it.

If X has a binomial distribution with n trials and probability of success p on each trial, then:

1. The mean of X is $\mu = np$.
2. The variance of X is $\sigma^2 = np(1-p)$.
3. The standard deviation of X is $\sigma = \sqrt{np(1-p)}$.

For example, suppose you flip a fair coin 100 times and let X be the number of heads; then X has a binomial distribution with $n = 100$ and $p = 0.50$. Its mean is $\mu = np = 100(0.50) = 50$ heads (which makes sense, because heads and tails are 50-50). The variance of X is $\sigma^2 = np(1-p) = 100(0.50)(1-0.50) = 25$, which is in square units (so you can't interpret it); and the standard deviation is the square root of the variance, which is 5. That means when you flip a coin 100 times, and do that over and over, the average number of heads you'll get is 50, and you can expect that to vary by about 5 heads on average.



The formula for the mean of a binomial distribution has intuitive meaning. The p in the formula represents the probability of a success, yes, but it also represents the *proportion* of successes you can expect in n trials. Therefore, the total *number* of successes you can expect — that is, the mean of X — is $\mu = np$.

The formula for variance has intuitive meaning as well. The only variability in the outcomes of each trial is between success (with probability p) and failure (with probability $1 - p$). Over n trials, the variance of the number of successes/failures is measured by $\sigma^2 = np(1-p)$. The standard deviation is just the square root.



If the value of n is too large to use the binomial formula or the binomial table to calculate probabilities (see the earlier sections in this chapter), there's an alternative. It turns out that if n is large enough, you can use the normal distribution to get an approximate answer for a binomial probability. The mean and standard deviation of the binomial are involved in this process. All the details are in Chapter 9.

Chapter 9

The Normal Distribution

In This Chapter

- ▶ Understanding the normal and standard normal distributions
 - ▶ Going from start to finish when finding normal probabilities
 - ▶ Working backward to find percentiles
-

In your statistical travels you'll come across two major types of random variables: discrete and continuous. *Discrete random variables* basically count things (number of heads on 10 coin flips, number of female Democrats in a sample, and so on). The most well-known discrete random variable is the binomial. (See Chapter 8 for more on discrete random variables and binomials). A *continuous random variable* is typically based on measurements; it either takes on an uncountably infinite number of values (values within an interval on the real line), or it has so many possible values that it may as well be deemed continuous (for example, time to complete a task, exam scores, and so on).

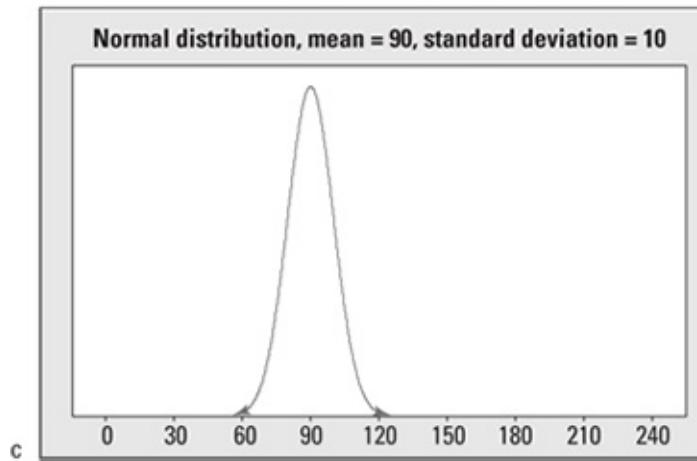
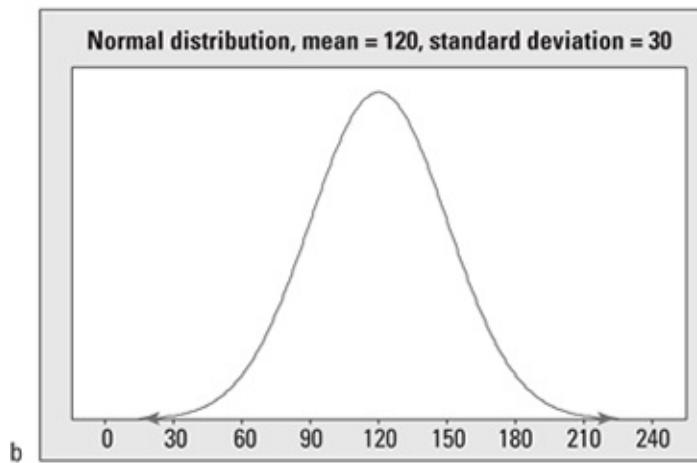
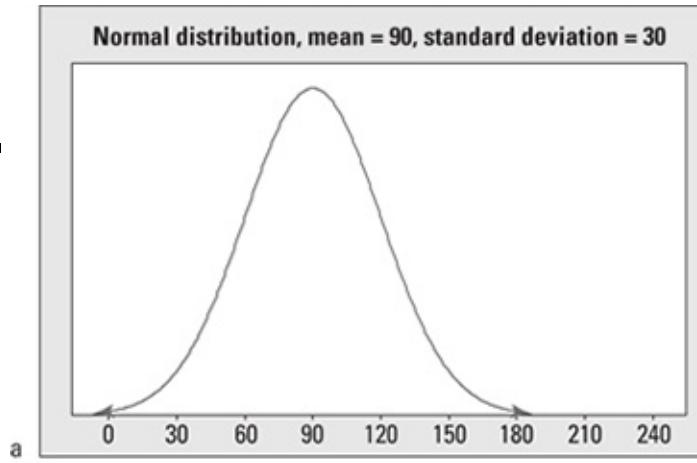
In this chapter, you understand and calculate probabilities for the most famous continuous random variable of all time — the normal distribution. You also find percentiles for the normal distribution, where you are given a probability as a percent and you have to find the value of X that's associated with it. And you can think how funny it would be to see a statistician wearing a T-shirt that said "I'd rather be normal."

Exploring the Basics of the Normal Distribution

A continuous random variable X has a normal distribution if its values fall into a smooth (continuous) curve with a bell-shaped pattern. Each normal distribution has its own mean, denoted by the Greek letter μ (say "mu"); and its own standard deviation, denoted by the Greek letter σ (say "sigma"). But no matter what their means and standard deviations are, all normal distributions have the same basic bell shape. Figure 9-1 shows some examples of normal distributions.

Figure 9-1:
Three normal distributions, with means and standard deviations of a) 90 and 30;

b) 120 and 30;
and c) 90 and
10,
respectively.



Every normal distribution has certain properties. You use these properties to determine the relative standing of any particular result on the distribution, and to find probabilities. The properties of any normal distribution are as follows:

- ✓ Its shape is symmetric (that is, when you cut it in half the two pieces are mirror images of each other).
- ✓ Its distribution has a bump in the middle, with tails going down and out to the left and right.
- ✓ The mean and the median are the same and lie directly in the middle of the distribution (due to symmetry).

- ✓ Its standard deviation measures the distance on the distribution from the mean to the *inflection point* (the place where the curve changes from an “upside-down-bowl” shape to a “right-side-up-bowl” shape).
- ✓ Because of its unique bell shape, probabilities for the normal distribution follow the Empirical Rule (full details in Chapter 5), which says the following:
 - About 68 percent of its values lie within one standard deviation of the mean. To find this range, take the value of the standard deviation, then find the mean plus this amount, and the mean minus this amount.
 - About 95 percent of its values lie within two standard deviations of the mean. (Here you take 2 times the standard deviation, then add it to and subtract it from the mean.)
 - Almost all of its values (about 99.7 percent of them) lie within three standard deviations of the mean. (Take 3 times the standard deviation and add it to and subtract it from the mean.)
- ✓ Precise probabilities for all possible intervals of values on the normal distribution (not just for those within 1, 2, or 3 standard deviations from the mean) are found using a table with minimal (if any) calculations. (The next section gives you all the info on this table.)

Take a look again at Figure 9-1. To compare and contrast the distributions shown in Figure 9-1a, b, and c, you first see they are all symmetric with the signature bell shape. The examples in Figure 9-1a and Figure 9-1b have the same standard deviation, but their means are different; Figure 9-1b is located 30 units to the right of Figure 9-1a because its mean is 120 compared to 90. Figures 9-1a and c have the same mean (90), but Figure 9-1a has more variability than Figure 9-1c due to its higher standard deviation (30 compared to 10). Because of the increased variability, the values in Figure 9-1a stretch from 0 to 180 (approximately), while the values in Figure 9-1c only go from 60 to 120.

Finally, Figures 9-1b and c have different means and different standard deviations entirely; Figure 9-1b has a higher mean which shifts it to the right, and Figure 9-1c has a smaller standard deviation; its values are the most concentrated around the mean.



Noting the mean and standard deviation is important so you can properly interpret numbers located on a particular normal distribution. For example, you can compare where the number 120 falls on each of the normal distributions in Figure 9-1. In Figure 9-1a, the number 120 is one standard deviation above the mean (because the standard deviation is 30, you get $90 + 1 * 30 = 120$). So on this first distribution, the number 120 is the upper value for the range where about 68% of the data are located, according to the Empirical Rule (see Chapter 5).

In Figure 9-1b, the number 120 lies directly on the mean, where the values are most

concentrated. In Figure 9-1c, the number 120 is way out on the rightmost fringe, 3 standard deviations above the mean (because the standard deviation this time is 10, you get $90 + 3[10] = 120$). In Figure 9-1c, values beyond 120 are very unlikely to occur because they are beyond the range where about 99.7% of the values should be, according to the Empirical Rule.

Meeting the Standard Normal (Z-) Distribution

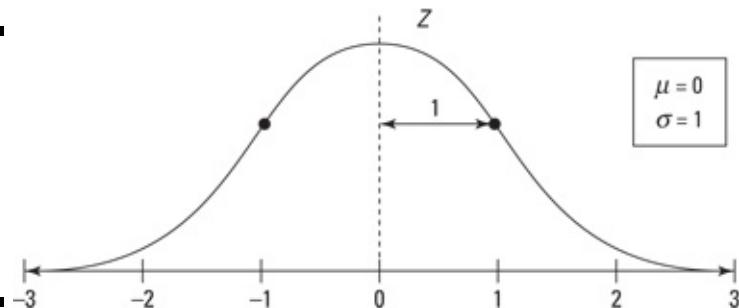
One very special member of the normal distribution family is called the standard normal distribution, or Z-distribution. The Z-distribution is used to help find probabilities and percentiles for regular normal distributions (X). It serves as the standard by which all other normal distributions are measured.

Checking out Z

The Z-distribution is a normal distribution with mean zero and standard deviation 1; its graph is shown in Figure 9-2. Almost all (about 99.7%) of its values lie between -3 and $+3$ according to the Empirical Rule. Values on the Z-distribution are called z -values, z -scores, or standard scores. A z -value represents the number of standard deviations that a particular value lies above or below the mean. For example, $z = 1$ on the Z-distribution represents a value that is 1 standard deviation above the mean. Similarly, $z = -1$ represents a value that is one standard deviation below the mean (indicated by the minus sign on the z -value). And a z -value of 0 is — you guessed it — right on the mean. All z -values are universally understood.

If you refer back to Figure 9-1 and the discussion regarding where the number 120 lies on each normal distribution in “Exploring the Basics of the Normal Distribution,” you can now calculate z -values to get a much clearer picture. In Figure 9-1a, the number 120 is located one standard deviation above the mean, so its z -value is 1. In Figure 9-1b, 120 is equal to the mean, so its z -value is 0. Figure 9-1c shows that 120 is 3 standard deviations above the mean, so its z -value is 3.

Figure 9-2:
The Z-distribution has a mean of 0 and standard deviation of 1.





High standard scores (z-values) aren't always the best. For example, if you're measuring the amount of time needed to run around the block, a standard score of +2 is a bad thing because your time was two standard deviations above (more than) the overall average time. In this case, a standard score of -2 would be much better, indicating your time was two standard deviations below (less than) the overall average time.

Standardizing from X to Z

Probabilities for any continuous distribution are found by finding the area under a curve (if you're into calculus, you know that means integration; if you're not into calculus, don't worry about it). Although the bell-shaped curve for the normal distribution looks easy to work with, calculating areas under its curve turns out to be a nightmare requiring high-level math procedures (believe me, I won't be going there in this book!). Plus, every normal distribution is different, causing you to repeat this process over and over each time you have to find a new probability.

To help get over this obstacle, statisticians worked out all the math gymnastics for one particular normal distribution, made a table of its probabilities, and told the rest of us to knock ourselves out. Can you guess which normal distribution they chose to crank out the table for?

Yes, all the basic results you need to find probabilities for any normal distribution (X) can be boiled down into one table based on the standard normal (Z -) distribution. This table is called the Z -table and is found in the appendix. Now all you need is one formula that transforms values from your normal distribution (X) to the Z -distribution; from there you can use the Z -table to find any probability you need.

Changing an x -value to a z -value is called *standardizing*. The so-called "z-formula" for standardizing an x -value to a z -value is:

$$z = \frac{x - \mu}{\sigma}$$

You take your x -value, subtract the mean of X , and divide by the standard deviation of X . This gives you the corresponding standard score (z -value or z -score).



Standardizing is just like changing units (for example, from Fahrenheit to Celsius). It doesn't affect probabilities for X ; that's why you can use the Z -table to find them!



You can standardize an x -value from any distribution (not just the normal) using the z -formula. Similarly, not all standard scores come from a normal distribution.



Because you subtract the mean from your x -values and divide everything by the standard deviation when you standardize, you are literally taking the mean and standard deviation of X out of the equation. This is what allows you to compare everything on the scale from -3 to $+3$ (the Z-distribution) where negative values indicate being below the mean, positive values indicate being above the mean, and a value of 0 indicates you're right on the mean.

Standardizing also allows you to compare numbers from different distributions. For example, suppose Bob scores 80 on both his math exam (which has a mean of 70 and standard deviation of 10) and his English exam (which has a mean of 85 and standard deviation of 5). On which exam did Bob do better, in terms of his relative standing in the class?

Bob's math exam score of 80 standardizes to a z -value of $\frac{80-70}{10} = \frac{10}{10} = 1$. That tells us his math score is one standard deviation above the class average. His English exam score of 80 standardizes to a z -value of $\frac{80-85}{5} = \frac{-5}{5} = -1$, putting him one standard deviation below the class average. Even though Bob scored 80 on both exams, he actually did better on the math exam than the English exam, relatively speaking.



To interpret a standard score, you don't need to know the original score, the mean, or the standard deviation. The standard score gives you the relative standing of a value, which in most cases is what matters most. In fact, on most national achievement tests, they won't even tell you what the mean and standard deviation were when they report your results; they just tell you where you stand on the distribution by giving you your z -score.

Finding probabilities for Z with the Z-table

A full set of less-than probabilities for a wide range of z -values is in the Z-table (Table A-1 in the appendix). To use the Z-table to find probabilities for the standard normal (Z-) distribution, do the following:

- 1. Go to the row that represents the first digit of your z -value and the first digit after the decimal point.**
- 2. Go to the column that represents the second digit after the decimal point of your z -value.**
- 3. Intersect the row and column.**

This result represents $p(Z < z)$, the probability that the random variable Z is less than the number z (also known as the percentage of z -values that are less than yours).

For example, suppose you want to find $p(Z < 2.13)$. Using the Z-table, find the row for 2.1

and the column for 0.03. Intersect that row and column to find the probability: 0.9834. You find that $p(Z < 2.13) = 0.9834$.

Suppose you want to look for $p(Z < -2.13)$. You find the row for -2.1 and the column for 0.03. Intersect the row and column and you find 0.0166; that means $p(Z < -2.13)$ equals 0.0166. (This happens to be one minus the probability that Z is less than 2.13 because $p(Z < +2.13)$ equals 0.9834. That's true because the normal distribution is symmetric; more on that in the following section.)

Finding Probabilities for a Normal Distribution

Here are the steps for finding a probability when X has any normal distribution:

1. **Draw a picture of the distribution.**
2. **Translate the problem into one of the following: $p(X < a)$, $p(X > b)$, or $p(a < X < b)$.**
3. **Shade in the area on your picture.**
4. **Standardize a (and/or b) to a z -score using the z -formula:**

$$z = \frac{x - \mu}{\sigma}$$

4. **Look up the z -score on the Z-table (Table A-1 in the appendix) and find its corresponding probability.**

(See the section “Standardizing from X to Z ” for more on the Z-table).

- 5a. **If you need a “less-than” probability — that is, $p(X < a)$ — you’re done.**
- 5b. **If you want a “greater-than” probability — that is, $p(X > b)$ — take one minus the result from Step 4.**
- 5c. **If you need a “between-two-values” probability — that is, $p(a < X < b)$ — do Steps 1–4 for b (the larger of the two values) and again for a (the smaller of the two values), and subtract the results.**



The probability that X is equal to any single value is 0 for any continuous random variable (like the normal). That's because continuous random variables consider probability as being area under the curve, and there's no area under a curve at one single point. This isn't true of discrete random variables.

Suppose, for example, that you enter a fishing contest. The contest takes place in a pond where the fish lengths have a normal distribution with mean $\mu = 16$ inches and standard deviation $\sigma = 4$ inches.

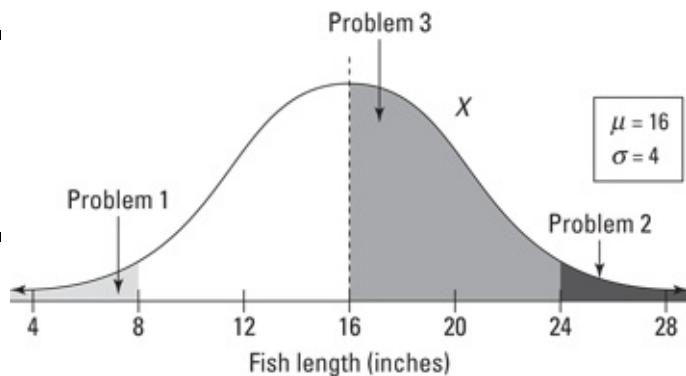
- ✓ Problem 1: What's the chance of catching a small fish — say, less than 8 inches?
- ✓ Problem 2: Suppose a prize is offered for any fish over 24 inches. What's the

chance of winning a prize?

✓ Problem 3: What's the chance of catching a fish between 16 and 24 inches?

To solve these problems using the steps that I just listed, first draw a picture of the normal distribution at hand. Figure 9-3 shows a picture of X 's distribution for fish lengths. You can see where the numbers of interest (8, 16, and 24) fall.

Figure 9-3:
The distribution of fish lengths in a pond.



Next, translate each problem into probability notation. Problem 1 is really asking you to find $p(X < 8)$. For Problem 2, you want $p(X > 24)$. And Problem 3 is looking for $p(16 < X < 24)$.

Step 3 says change the x -values to z -values using the z -formula:

$$z = \frac{x - \mu}{\sigma}$$

For Problem 1 of the fish example, you have the following:

$$p(X < 8) = p\left(Z < \frac{8-16}{4}\right) = p(Z < -2)$$

Similarly for Problem 2, $p(X > 24)$ becomes

$$p(X > 24) = p\left(Z > \frac{24-16}{4}\right) = p(Z > 2)$$

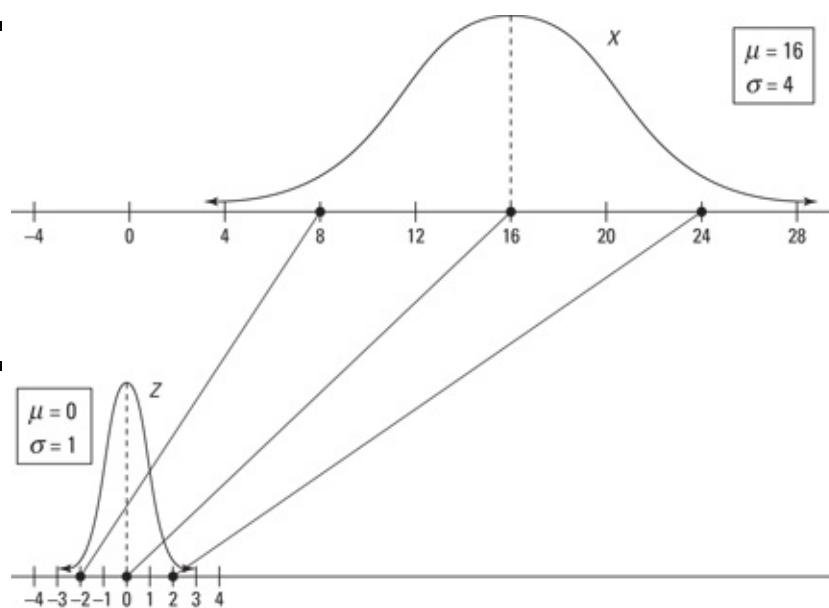
And Problem 3 translates from $p(16 < X < 24)$ to

$$p(16 < X < 24) = p\left(\frac{16-16}{4} < Z < \frac{24-16}{4}\right) = p(0 < Z < 2)$$

Figure 9-4 shows a comparison of the X -distribution and Z -distribution for the values $x = 8, 16$, and 24 , which standardize to $z = -2, 0$, and $+2$, respectively.

Now that you have changed x -values to z -values, you move to Step 4 and find (or calculate) probabilities for those z -values using the Z -table (in the appendix). In Problem 1 of the fish example, you want $p(Z < -2)$; go to the Z -table and look at the row for -2.0 and the column for 0.00 , intersect them, and you find 0.0228 — according to Step 5a, you're done. The chance of a fish being less than 8 inches is equal to 0.0228 .

Figure 9-4:
Standardizing
numbers
from a normal
distribution
(X) to
numbers on
the Z -
distribution.



For Problem 2, find $p(Z > 2.00)$. Because it's a “greater-than” problem, this calls for Step 5b. To be able to use the Z -table, you need to rewrite this in terms of a “less-than” statement. Because the entire probability for the Z -distribution equals 1, we know $p(Z > 2.00) = 1 - p(Z < 2.00) = 1 - 0.9772 = 0.0228$ (using the Z -table). So, the chance that a fish is greater than 24 inches is also 0.0228. (Note: The answers to Problems 1 and 2 are the same because the Z -distribution is symmetric; refer to Figure 9-3.)

In Problem 3, you find $p(0 < Z < 2.00)$; this requires Step 5c. First find $p(Z < 2.00)$, which is 0.9772 from the Z -table. Then find $p(Z < 0)$, which is 0.5000 from the Z -table. Subtract them to get $0.9772 - 0.5000 = 0.4772$. The chance of a fish being between 16 and 24 inches is 0.4772.



The Z -table does not list every possible value of Z ; it just carries them out to two digits after the decimal point. Use the one closest to the one you need. And just like in an airplane where the closest exit may be behind you, the closest z -value may be the one that is lower than the one you need.

Finding X When You Know the Percent

Another popular normal distribution problem involves finding percentiles for X (see Chapter 5 for a detailed rundown on percentiles). That is, you are given the percentage or probability of being at or below a certain x -value, and you have to find the x -value that corresponds to it. For example, if you know that the people whose golf scores were in the lowest 10% got to go to the tournament, you may wonder what the cutoff score was; that score would represent the 10th percentile.



A percentile isn't a percent. A percent is a number between 0 and 100; a percentile

is a value of X (a height, an IQ, a test score, and so on).

Figuring out a percentile for a normal distribution

Certain percentiles are so popular that they have their own names and their own notation. The three “named” percentiles are Q_1 — the first quartile, or the 25th percentile; Q_2 — the 2nd quartile (also known as the *median* or the 50th percentile); and Q_3 — the 3rd quartile or the 75th percentile. (See Chapter 5 for more information on quartiles.)

Here are the steps for finding any percentile for a normal distribution X :

- 1a. If you’re given the probability (percent) less than x and you need to find x , you translate this as: Find a where $p(X < a) = p$ (and p is the given probability). That is, find the p^{th} percentile for X . Go to Step 2.
- 1b. If you’re given the probability (percent) greater than x and you need to find x , you translate this as: Find b where $p(X > b) = p$ (and p is given). Rewrite this as a percentile (less-than) problem: Find b where $p(X < b) = 1 - p$. This means find the $(1 - p)^{\text{th}}$ percentile for X .
2. Find the corresponding percentile for Z by looking in the body of the Z -table (in the appendix) and finding the probability that is closest to p (from Step 1a) or $1 - p$ (from Step 1b). Find the row and column this probability is in (using the table backwards). This is the desired z -value.
3. Change the z -value back into an x -value (original units) by using $x = \mu + z\sigma$. You’ve (finally!) found the desired percentile for X .

The formula in this step is just a rewriting of the z -formula, $z = \frac{x - \mu}{\sigma}$, so it’s solved for x .

Doing a low percentile problem

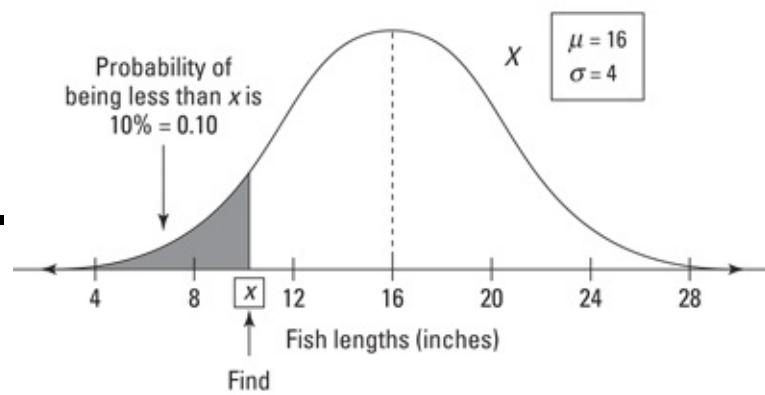
Look at the fish example used previously in “Finding Probabilities for a Normal Distribution,” where the lengths (X) of fish in a pond have a normal distribution with mean 16 inches and standard deviation 4 inches. Suppose you want to know what length marks the bottom 10 percent of all the fish lengths in the pond. What percentile are you looking for?



Being at the bottom 10 percent means you have a “less-than” probability that’s equal to 10 percent, and you are at the 10th percentile.

Now go to Step 1a in the preceding section and translate the problem. In this case, because you’re dealing with a “less-than” situation, you want to find x such that $p(X < x) = 0.10$. This represents the 10th percentile for X . Figure 9-5 shows a picture of this situation.

Figure 9-5:
Bottom 10 percent of fish in the pond, according to length.



Now go to Step 2, which says to find the 10th percentile for Z . Looking in the body of the Z -table (in the appendix), the probability closest to 0.10 is 0.1003, which falls in the row for $z = -1.2$ and the column for 0.08. That means the 10th percentile for Z is -1.28 ; so a fish whose length is 1.28 standard deviations below the mean marks the bottom 10 percent of all fish lengths in the pond.

But exactly how long is that fish, in inches? In Step 3, you change the z -value back to an x -value (fish length in inches) using the z -formula solved for x ; you get $x = 16 + -1.28 * 4 = 10.88$ inches. So 10.88 inches marks the lowest 10 percent of fish lengths. Ten percent of the fish are shorter than that.

Working with a higher percentile

Now suppose you want to find the length that marks the *top* 25 percent of all the fish in the pond. This problem calls for Step 1b (in “Finding a percentile for a normal distribution”) because being in the top part of the distribution means you’re dealing with a greater-than probability. The number you are looking for is somewhere in the right tail (upper area) of the X -distribution, with $p = 25$ percent of the probability to its right and $1 - p = 75$ percent to its left. Thinking in terms of the Z -table and how it only uses less-than probabilities, you need to find the 75th percentile for Z , then change it to an x -value.

Step 2: The 75th percentile of Z is the z -value where $p(Z < z) = 0.75$. Using the Z -table (in the appendix), you find the probability closest to 0.7500 is 0.7486, and its corresponding z -value is in the row for 0.6 and column for 0.07. Put these together and you get a z -value of 0.67. This is the 75th percentile for Z . In Step 3, change the z -value back to an x -value (length in inches) using the z -formula solved for x to get $x = 16 + 0.67 * 4 = 18.68$ inches. So, 75% of the fish are shorter than 18.68 inches. And to answer the original question, the top 25% of the fish in the pond are longer than 18.68 inches.

Translating tricky wording in percentile problems



Some percentile problems are especially challenging to translate. For example,

Suppose the amount of time for a racehorse to run around a track in a qualifying round has a normal distribution with mean 120 seconds and standard deviation 5 seconds. The best 10 percent of the times qualify; the rest don't. What's the cutoff time for qualifying?

Because “best times” mean “lowest times” in this case, the percentage of times that lie *below* the cutoff must be 10, and the percentage *above* the cutoff must be 90. (It’s an easy mistake to think it’s the other way around.) The percentile of interest is therefore the 10th, which is down on the left tail of the distribution. You now work this problem the same way I worked Problem 1 regarding fish lengths (see the section, “[Finding Probabilities for a Normal Distribution](#)”). The standard score for the 10th percentile is $z = -1.28$ looking at the Z-table (in the appendix). Converting back to original units, you get $x = \mu + z\sigma = 120 + (-1.28)(5) = 113.6$ seconds. So the cutoff time needed for a racehorse to qualify (that is, to be among the fastest 10%) is 113.6 seconds. (Notice this number is less than the average time of 120 seconds, which makes sense; a negative z -value is what makes this happen.)



The 50th percentile for the normal distribution is the mean (because of symmetry) and its z -score is zero. Smaller percentiles, like the 10th, lie below the mean and have negative z -scores. Larger percentiles, like the 75th, lie above the mean and have positive z -scores.

Here’s another style of wording that has a bit of a twist: Suppose times to complete a statistics exam have a normal distribution with a mean of 40 minutes and standard deviation of 6 minutes. Deshawn’s time comes in at the 90th percentile. What percentage of the students are still working on their exams when Deshawn leaves? Because Deshawn is at the 90th percentile, 90 percent of the students have exam times lower than hers. That means 90% of the students left before Deshawn, so $100 - 90 = 10$ percent of the students are still working when Deshawn leaves.



To be able to decipher the language used to imply a percentile problem, look for clues like *the bottom 10%* (also known as the 10th percentile) and *the top 10%* (also known as the 90th percentile). For *the best 10%*, you must determine whether low or high numbers qualify as “best.”

Normal Approximation to the Binomial

Suppose you flip a fair coin 100 times and you let X equal the number of heads. What’s the probability that X is greater than 60? In Chapter 8, you solve problems like this (involving fewer flips) using the binomial distribution. For binomial problems where n (the number of trials) is small, you can either use the direct formula (found in Chapter 8),

the binomial table (found in the appendix), or you can use technology if available (such as a graphing calculator or Microsoft Excel).

However, if n is large the calculations get unwieldy and the binomial table runs out of numbers. If there's no technology available (like when taking an exam), what can you do to find a binomial probability? Turns out, if n is large enough, you can use the normal distribution to find a very close approximate answer with a lot less work.

But what do I mean by n being “large enough”? To determine whether n is large enough to use what statisticians call the *normal approximation to the binomial*, both of the following conditions must hold:

- ✓ $n * p \geq 10$ (at least 10), where p is the probability of success
- ✓ $n * (1 - p) \geq 10$ (at least 10), where $1 - p$ is the probability of failure

To find the normal approximation to the binomial distribution when n is large, use the following steps:

1. Verify whether n is large enough to use the normal approximation by checking the two appropriate conditions.

For the coin-flipping question, the conditions are met because $n * p = 100 * 0.50 = 50$, and $n * (1 - p) = 100 * (1 - 0.50) = 50$, both of which are at least 10. So go ahead with the normal approximation.

2. Translate the problem into a probability statement about X .

For the coin-flipping example, you need to find $p(X > 60)$.

3. Standardize the x -value to a z -value, using the z -formula:

$$z = \frac{x - \mu}{\sigma}$$

For the mean of the normal distribution, use $\mu = np$ (the mean of the binomial), and for the standard deviation σ , use $\sqrt{np(1-p)}$ (the standard deviation of the binomial; see Chapter 8).

For the coin-flipping example, use $\mu = np = (100)(0.50) = 50$ and $\sigma = \sqrt{np(1-p)} = \sqrt{100(0.50)(1-0.50)} = 5$.

Then put these values into the z -formula to get $z = \frac{x - \mu}{\sigma} = \frac{60 - 50}{5} = 2$. To solve the problem, you need to find $p(Z > 2)$.



On an exam, you won't see μ and σ in the problem when you have a binomial distribution. However, you know the formulas that allow you to calculate both of them using n and p (both of which will be given in the problem). Just remember you have to do that extra step to calculate the μ and σ needed for the z -formula.

4. Proceed as you usually would for any normal distribution. That is, do Steps 4 and 5 described in the earlier section “Finding Probabilities for a Normal Distribution.”

Continuing the example, $p(Z > 2.00) = 1 - 0.9772 = 0.0228$ from the Z -table (appendix).

So the chance of getting more than 60 heads in 100 flips of a coin is only about 2.28

percent. (I wouldn't bet on it.)



When using the normal approximation to find a binomial probability, your answer is an *approximation* (not exact) — be sure to state that. Also show that you checked both necessary conditions for using the normal approximation.

Chapter 10

The *t*-Distribution

In This Chapter

- ▶ Characteristics of the *t*-distribution
 - ▶ Relationship between *Z*- and *t*-distributions
 - ▶ Understanding and using the *t*-table
-

The *t*-distribution is one of the mainstays of data analysis. You may have heard of the “*t*-test” for example, which is often used to compare two groups in medical studies and scientific experiments.

This short chapter covers the basic characteristics and uses of the *t*-distribution. You find out how it compares to the normal distribution (more on that in Chapter 9) and how to use the *t*-table to find probabilities and percentiles.

Basics of the *t*-Distribution

In this section, you get an overview of the *t*-distribution, its main characteristics, when it’s used, and how it’s related to the *Z*-distribution (see Chapter 9).

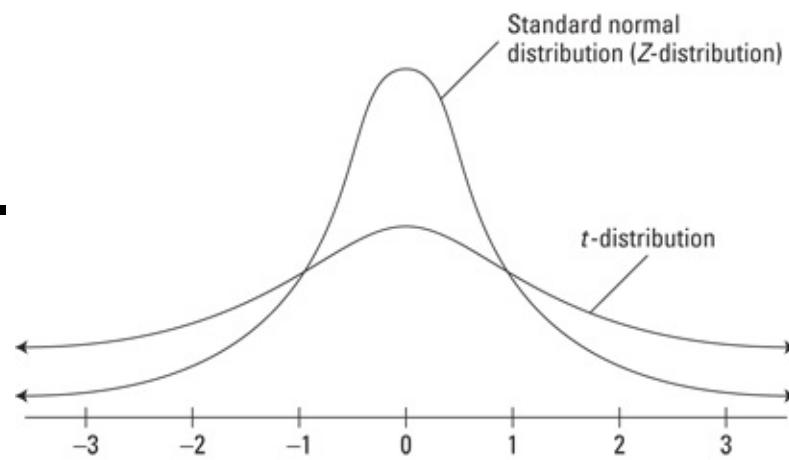
Comparing the *t*- and *Z*-distributions

The normal distribution is that well-known bell-shaped distribution whose mean is μ and whose standard deviation is σ (see Chapter 9 for more on the normal distribution). The most common normal distribution is the standard normal (also called the *Z*-distribution), whose mean is 0 and standard deviation is 1.

The *t*-distribution can be thought of as a cousin of the standard normal distribution — it looks similar in that it’s centered at zero and has a basic bell-shape, but it’s shorter and flatter than the *Z*-distribution. Its standard deviation is proportionally larger compared to the *Z*, which is why you see the fatter tails on each side.

Figure 10-1 compares the *t*- and standard normal (*Z*) distributions in their most general forms.

Comparing
the standard
normal (Z-)
distribution to
a generic *t*-
distribution.



The *t*-distribution is typically used to study the mean of a population, rather than to study the individuals within a population. In particular, it is used in many cases when you use data to estimate the population mean — for example, to estimate the average price of all the new homes in California. Or when you use data to test someone's claim about the population mean — for example, is it true that the mean price of all the new homes in California is \$500,000?



These procedures are called *confidence intervals* and *hypothesis tests* and are discussed in Chapters 13 and 14, respectively.

The connection between the normal distribution and the *t*-distribution is that the *t*-distribution is often used for analyzing the mean of a population if the population has a normal distribution (or fairly close to it). Its role is especially important if your data set is small or if you don't know the standard deviation of the population (which is often the case).

When statisticians use the term *t-distribution*, they aren't talking about just one individual distribution. There is an entire family of specific *t*-distributions, depending on what sample size is being used to study the population mean. Each *t*-distribution is distinguished by what statisticians call its *degrees of freedom*. In situations where you have one population and your sample size is n , the degrees of freedom for the corresponding *t*-distribution is $n - 1$. For example, a sample of size 10 uses a *t*-distribution with $10 - 1$, or 9, degrees of freedom, denoted t_9 (pronounced *tee sub-nine*). Situations involving two populations use different degrees of freedom and are discussed in Chapter 15.

Discovering the effect of variability on *t*-distributions

t-distributions based on smaller sample sizes have larger standard deviations than those based on larger sample sizes. Their shapes are flatter; their values are more spread out. That's because results based on smaller data sets are more variable than results based

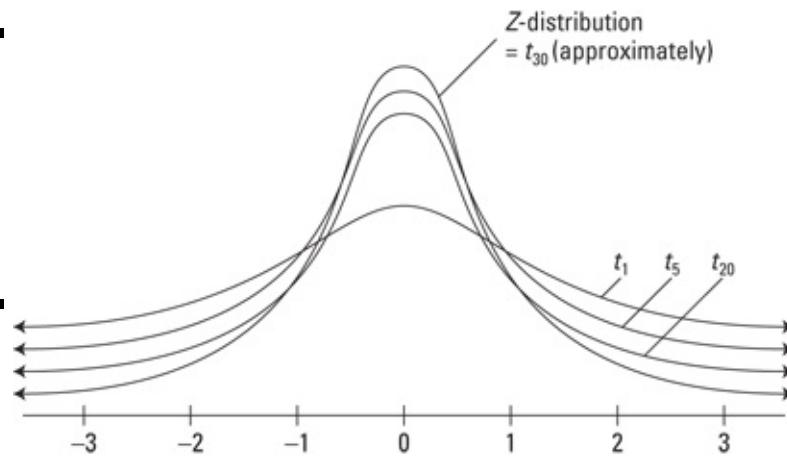
on large data sets.



The larger the sample size is, the larger the degrees of freedom will be, and the more the *t*-distributions look like the standard normal distribution (*Z*-distribution). A rough cutoff point where the *t*- and *Z*-distributions become similar (at least similar enough for jazz or government work) is around $n = 30$.

Figure 10-2 shows what different *t*-distributions look like for different sample sizes and how they all compare to the standard normal (*Z*) distribution.

Figure 10-2: *t*-distributions for different sample sizes compared to the *Z*-distribution.



Using the *t*-Table

Each normal distribution has its own mean and standard deviation that classify it, so finding probabilities for each normal distribution on its own is not the way to go.

Thankfully, you can standardize the values of any normal distribution to become values on a standard normal (*Z*) distribution (whose mean is 0 and standard deviation is 1) and use a *Z*-table (in the appendix) to find probabilities. (Chapter 9 has info on normal distributions.)

In contrast, a *t*-distribution is not classified by its mean and standard deviation, but by the sample size of the data set being used (n). Unfortunately, there is no single “standard *t*-distribution” that you can use to transform the numbers and find probabilities on a table. Because it wouldn’t be humanly possible to create a table of probabilities and corresponding *t*-values for every possible *t*-distribution, statisticians created one table showing certain values of *t*-distributions for a selection of degrees of freedom and a selection of probabilities. This table is called the *t*-table (it appears in the appendix). In this section, you find out how to find probabilities, percentiles, and critical values (for confidence intervals) using the *t*-table.

Finding probabilities with the *t*-table

Each row of the t -table (in the appendix) represents a different t -distribution, classified by its degrees of freedom (df). The columns represent various common greater-than probabilities, such as 0.40, 0.25, 0.10, and 0.05. The numbers across a row indicate the values on the t -distribution (the t -values) corresponding to the greater-than probabilities shown at the top of the columns. Rows are arranged by degrees of freedom.



Another term for greater-than probability is *right-tail probability*, which indicates that such probabilities represent areas on the right-most end (tail) of the t -distribution.

For example, the second row of the t -table is for the t_2 distribution (2 degrees of freedom, pronounced *tee sub-two*). You see that the second number, 0.816, is the value on the t_2 distribution whose area to its right (its right-tail probability) is 0.25 (see the heading for column 2). In other words, the probability that t_2 is greater than 0.816 equals 0.25. In probability notation, that means $p(t_2 > 0.816) = 0.25$.

The next number in row two of the t -table is 1.886, which lies in the 0.10 column. This means the probability of being greater than 1.886 on the t_2 distribution is 0.10. Because 1.886 falls to the right of 0.816, its right-tail probability is lower.

Figuring percentiles for the t -distribution

You can also use the t -table (in the appendix) to find percentiles for a t -distribution. A *percentile* is a number on a distribution whose less-than probability is the given percentage; for example, the 95th percentile of the t -distribution with $n - 1$ degrees of freedom is that value of t_{n-1} whose left-tail (less-than) probability is 0.95 (and whose right-tail probability is 0.05). (See Chapter 5 for particulars on percentiles.)

Suppose you have a sample of size 10 and you want to find the 95th percentile of its corresponding t -distribution. You have $n - 1 = 9$ degrees of freedom, so you look at the row for $df = 9$. The 95th percentile is the number where 95% of the values lie below it and 5% lie above it, so you want the right-tail area to be 0.05. Move across the row, find the column for 0.05, and you get $t_9 = 1.833$. This is the 95th percentile of the t -distribution with 9 degrees of freedom.

Now, if you increase the sample size to $n = 20$, the value of the 95th percentile decreases; look at the row for $20 - 1 = 19$ degrees of freedom, and in the column for 0.05 (a right-tail probability of 0.05) you find $t_{19} = 1.729$. Notice that the 95th percentile for the t_{19} distribution is less than the 95th percentile for the t_9 distribution (1.833). This is because larger degrees of freedom indicate a smaller standard deviation and the t -values are more concentrated about the mean, so you reach the 95th percentile with a smaller value of t . (See the section “Discovering the effect of variability on t -distributions,” earlier in this chapter.)

Picking out t^* -values for confidence intervals

Confidence intervals estimate population parameters, such as the population mean, by using a statistic (for example, the sample mean) plus or minus a margin of error. (See Chapter 13 for all the information you need on confidence intervals and more.) To compute the margin of error for a confidence interval, you need a *critical value* (the number of standard errors you add and subtract to get the margin of error you want; see Chapter 13). When the sample size is large (at least 30), you use critical values on the Z -distribution (shown in Chapter 13) to build the margin of error. When the sample size is small (less than 30) and/or the population standard deviation is unknown, you use the t -distribution to find critical values.

To help you find critical values for the t -distribution, you can use the last row of the t -table, which lists common confidence levels, such as 80%, 90%, and 95%. To find a critical value, look up your confidence level in the bottom row of the table; this tells you which column of the t -table you need. Intersect this column with the row for your df (see Chapter 13 for degrees of freedom formulas). The number you see is the critical value (or the t^* -value) for your confidence interval. For example, if you want a t^* -value for a 90% confidence interval when you have 9 degrees of freedom, go to the bottom of the table, find the column for 90%, and intersect it with the row for $df = 9$. This gives you a t^* -value of 1.833 (rounded).



Across the top row of the t -table, you see right-tail probabilities for the t -distribution. But confidence intervals involve both left- and right-tail probabilities (because you add and subtract the margin of error). So half of the probability left from the confidence interval goes into each tail. You need to take that into account. For example, a t^* -value for a 90% confidence interval has 5% for its greater-than probability and 5% for its less-than probability (taking 100% minus 90% and dividing by 2). Using the top row of the t -table, you would have to look for 0.05 (rather than 10%, as you might be inclined to do.) But using the bottom row of the table, you just look for 90%. (The result you get using either method ends up being in the same column.)



When looking for t^* -values for confidence intervals, use the bottom row of the t -table as your guide, rather than the headings at the top of the table.

Studying Behavior Using the t-Table

You can use computer software to calculate any probabilities, percentiles, or critical values you need for any t -distribution (or any other distribution) if it's available to you.

(On exams it may not be available.) However, one of the nice things about using a table to find probabilities (rather than using computer software) is that the table can tell you information about the behavior of the distribution itself — that is, it can give you the big picture. Here are some nuggets of big-picture information about the *t*-distribution you can glean by scanning the *t*-table (in the appendix).

In Figure 10-2, as the degrees of freedom increase, the values on each *t*-distribution become more concentrated around the mean, eventually resembling the *Z*-distribution. The *t*-table confirms this pattern as well. Because of the way the *t*-table is set up, if you choose any column and move down through the numbers in the column, you're increasing the degrees of freedom (and sample size) and keeping the right-tail probability the same. As you do this, you see the *t*-values getting smaller and smaller, indicating the *t*-values are becoming closer to (hence more concentrated around) the mean.

I labeled the second-to-last row of the *t*-table with a *z* in the *df* column. This indicates the “limit” of the *t*-values as the sample size (*n*) goes to infinity. The *t*-values in this row are approximately the same as the *z*-values on the *Z*-table (in the appendix) that correspond to the same greater-than probabilities. This confirms what you already know: As the sample size increases, the *t*- and the *Z*-distributions look more and more alike. For example, the *t*-value in row 30 of the *t*-table corresponding to a right-tail probability of 0.05 (column 0.05) is 1.697. This lies close to *z* = 1.645, the value corresponding to a right-tail area of 0.05 on the *Z*-distribution. (See row Z of the *t*-table.)



It doesn't take a super-large sample size for the values on the *t*-distribution to get close to the values on a *Z*-distribution. For example, when *n* = 31 and *df* = 30, the values in the *t*-table are already quite close to the corresponding values on the *Z*-table.

Sampling Distributions and the Central Limit Theorem

In This Chapter

- ▶ Understanding the concept of a sampling distribution
 - ▶ Putting the Central Limit Theorem to work
 - ▶ Determining the factors that affect precision
-

When you take a sample of data, it's important to realize the results will vary from sample to sample. Statistical results based on samples should include a measure of how much those results are expected to vary. When the media reports statistics like the average price of a gallon of gas in the U.S. or the percentage of homes on the market that were sold over the last month, you know they didn't sample every possible gas station or every possible home sold. The question is, how much would their results change if another sample was selected?

This chapter addresses this question by studying the behavior of means for all possible samples, and the behavior of proportions from all possible samples. By studying the behavior of all possible samples, you can gauge where your sample results fall and understand what it means when your sample results fall outside of certain expectations.

Defining a Sampling Distribution

A *random variable* is a characteristic of interest that takes on certain values in a random manner. For example, the number of red lights you hit on the way to work or school is a random variable; the number of children a randomly selected family has is a random variable. You use capital letters such as X or Y to denote random variables and you use small case letters x or y to denote actual outcomes of random variables. A *distribution* is a listing, graph, or function of all possible outcomes of a random variable (such as X) and how often each actual outcome (x), or set of outcomes, occurs. (See Chapter 8 for more details on random variables and distributions.)

For example, suppose a million of your closest friends each rolls a single die and records each actual outcome (x). A table or graph of all these possible outcomes (one through six) and how often they occurred represents the distribution of the random variable X . A graph of the distribution of X in this case is shown in Figure 11-1a. It shows the numbers

1–6 appearing with equal frequency (each one occurring 1/6 of the time), which is what you expect over many rolls if the die is fair.

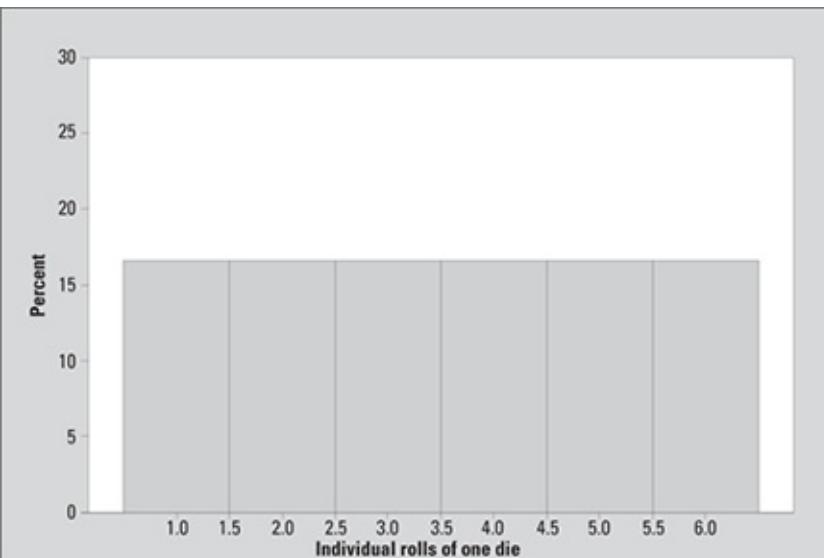
Now suppose each of your friends rolls this single die 50 times ($n = 50$) and records the average, \bar{x} . The graph of all their averages of all their samples represents the distribution of the random variable \bar{X} . Because this distribution is based on sample averages rather than individual outcomes, this distribution has a special name. It's called the *sampling distribution* of the sample mean, \bar{X} . Figure 11-1b shows the sampling distribution of \bar{X} , the average of 50 rolls of a die.

Figure 11-1b (average of 50 rolls) shows the same range (1 through 6) of outcomes as Figure 11-1a (individual rolls), but Figure 11-1b has more possible outcomes. You could get an average of 3.3 or 2.8 or 3.9 for 50 rolls, for example, whereas someone rolling a single die can only get whole numbers from 1 to 6. Also, the shape of the graphs are different; Figure 11-1a shows a flat shape, where each outcome is equally likely, and Figure 11-1b has a mound shape; that is, outcomes near the center (3.5) occur with high frequency and outcomes near the edges (1 and 6) occur with extremely low frequency. A detailed look at the differences and similarities in shape, center, and spread for individuals versus averages, and the reasons behind them, is the topic of the following sections. (See Chapter 8 if you need background info on shape, center, and spread of random variables before diving in.)

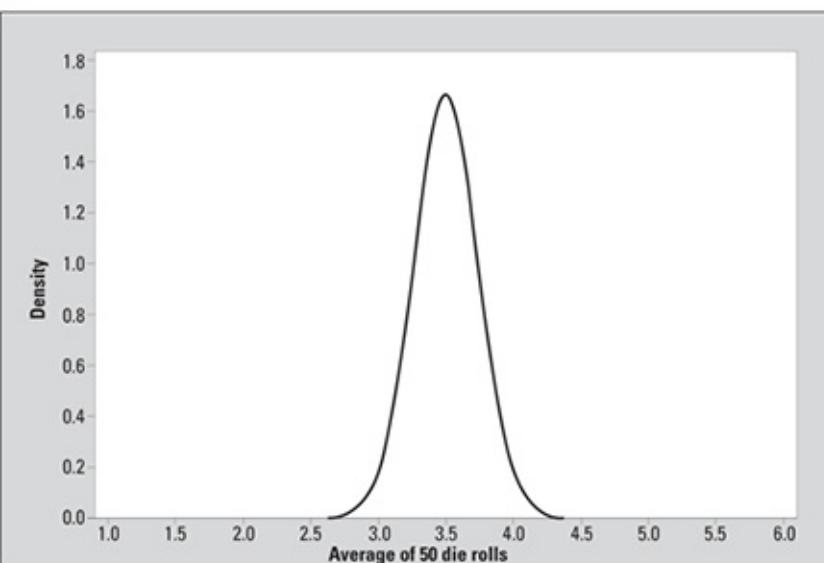
The Mean of a Sampling Distribution

Using the die-rolling example from the preceding section, X is a random variable denoting the outcome you can get from a single die (assuming the die is fair). The mean of X (over all possible outcomes) is denoted by μ_x (pronounced *mu sub-x*); in this case its value is 3.5 (as shown in Figure 11-1a). If you roll a die 50 times and take the average, the random variable \bar{X} represents any outcome you could get. The mean of \bar{X} , denoted $\mu_{\bar{x}}$ (pronounced *mu sub-x-bar*) equals 3.5 as well. (You can see this result in Figure 11-1b.)

Figure 11-1:
Distributions
of a)
individual
rolls of one
die; and b)
average of 50
rolls of one
die.



a



b

This result is no coincidence! In general, the mean of the population of all possible sample means is the same as the mean of the original population. (Notationally speaking, you write $\mu_{\bar{x}} = \mu_x$.) It's a mouthful, but it makes sense that the average of the averages from all possible samples is the same as the average of the population that the samples came from. In the die rolling example, the average of the population of all 50-roll averages equals the average of the population of all single rolls (3.5).



Using subscripts on μ , you can distinguish which mean you're talking about — the mean of X (all individuals in a population) or the mean of \bar{X} (all sample means from the population).

Measuring Standard Error

The values in any population deviate from their mean; for instance, people's heights differ from the overall average height. Variability in a population of individuals (X) is

measured in *standard deviations* (see Chapter 5 for details on standard deviation). Sample means vary because you're not sampling the whole population, only a subset; and as samples vary, so will their means. Variability in the sample mean (\bar{X}) is measured in terms of *standard errors*.



Error here doesn't mean there's been a mistake — it means there is a gap between the population and sample results.

The standard error of the sample mean is denoted by $\sigma_{\bar{x}}$ (*sigma sub-x-bar*). Its formula is $\frac{\sigma_x}{\sqrt{n}}$, where σ_x is population standard deviation (*sigma sub-x*) and n is size of each sample. In the next sections you see the effect each of these two components has on the standard error.

Sample size and standard error

The first component of standard error is the sample size, n . Because n is in the denominator of the standard error formula, the standard error decreases as n increases. It makes sense that having more data gives less variation (and more precision) in your results.

Suppose X is the time it takes for a clerical worker to type and send one letter of recommendation, and say X has a normal distribution with mean 10.5 minutes and standard deviation 3 minutes. The bottom curve in Figure 11-2 shows the picture of the distribution of X , the individual times for all clerical workers in the population. According to the Empirical Rule (see Chapter 9), most of the values are within 3 standard deviations of the mean (10.5) — between 1.5 and 19.5.

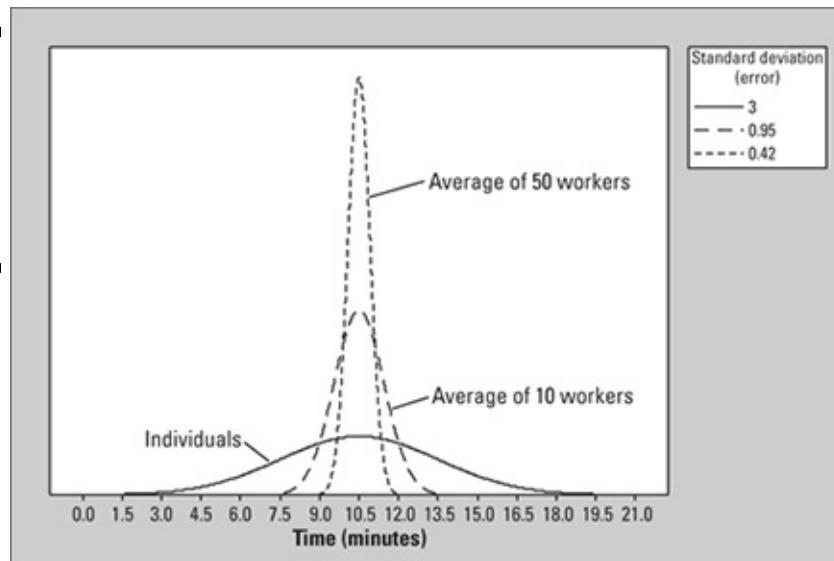
Now take a random sample of 10 clerical workers, measure their times, and find the average, \bar{x} , each time. Repeat this process over and over, and graph all the possible results for all possible samples. The middle curve in Figure 11-2 shows the picture of the sampling distribution of \bar{X} . Notice that it's still centered at 10.5 (which you expected) but its variability is smaller; the standard error in this case is $\frac{\sigma_x}{\sqrt{n}} = \frac{3}{\sqrt{10}} = 0.95$ minutes (quite a bit less than 3 minutes, the standard deviation of the individual times).

Looking at Figure 11-2, the average times for samples of 10 clerical workers are closer to the mean (10.5) than the individual times are. That's because average times don't change as much from sample to sample as individual times change from person to person.

Now take all possible random samples of 50 clerical workers and find their means; the sampling distribution is shown in the tallest curve in Figure 11-2. The standard error of \bar{X} goes down to $\frac{\sigma_x}{\sqrt{n}} = \frac{3}{\sqrt{50}} = 0.42$ minutes. You can see the average times for 50 clerical workers are even closer to 10.5 than the ones for 10 clerical workers. By the Empirical Rule, most

of the values fall between $10.5 - 3(.42) = 9.24$ and $10.5 + 3(.42) = 11.76$. Larger samples give even more precision around the mean because they change even less from sample to sample.

Figure 11-2:
Distributions
of times for 1
worker, 10
workers, and
50 workers.



Why is having more precision around the mean important? Because sometimes you don't know the mean but want to determine what it is, or at least get as close to it as possible. How can you do that? By taking a large random sample from the population and finding its mean. You know that your sample mean will be close to the actual population mean if your sample is large, as Figure 11-2 shows (assuming your data are collected correctly; see Chapter 16 for details on collecting good data).

Population standard deviation and standard error

The second component of standard error involves the amount of diversity in the population (measured by standard deviation). In the standard error formula $\frac{\sigma_x}{\sqrt{n}}$, for \bar{X} , you see the population standard deviation, σ_x , is in the numerator. That means as the population standard deviation increases, the standard error of the sample means also increases. Mathematically this makes sense; how about statistically?

Suppose you have two ponds full of fish (call them pond #1 and pond #2), and you're interested in the length of the fish in each pond. Assume the fish lengths in each pond have a normal distribution (see Chapter 9). You've been told that the fish lengths in pond #1 have a mean of 20 inches and a standard deviation of 2 inches (see Figure 11-3a). Suppose the fish in pond #2 also average 20 inches but have a larger standard deviation of 5 inches (see Figure 11-3b).

Comparing Figures 11-3a and 11-3b, you see the lengths for the two populations of fish have the same shape and mean, but the distribution in Figure 11-3b (for pond #2) has more spread, or variability, than the distribution shown in Figure 11-3a (for pond #1).

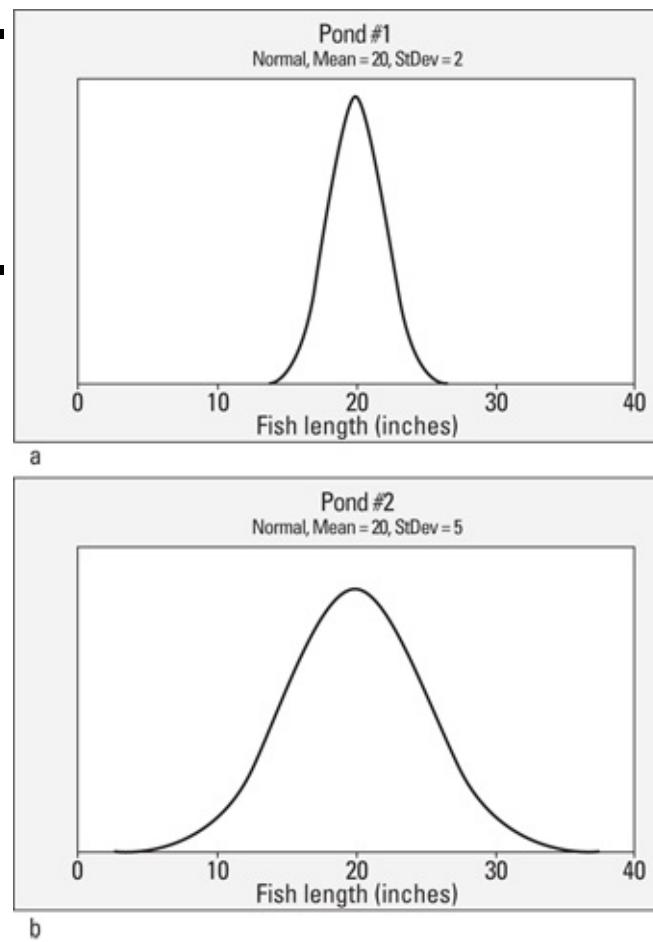
This spread confirms that the fish in pond #2 vary more in length than those in pond #1.

Now suppose you take a random sample of 100 fish from pond #1, find the mean length of the fish, and repeat this process over and over. Then you do the same with pond #2. Because the lengths of individual fish in pond #2 have more variability than the lengths of individual fish in pond #1, you know the average lengths of samples from pond #2 will have more variability than the average lengths of samples from pond #1 as well. (In fact, you can calculate their standard errors using the formula earlier in this section to be 0.20 and 0.50, respectively.)



Estimating the population average is harder when the population varies a lot to begin with — estimating the population average is much easier when the population values are more consistent. The bottom line is the standard error of the sample mean is larger when the population standard deviation is larger.

Figure 11-3:
Distributions
of fish
lengths a) in
pond #1; b) in
pond #2.



Looking at the Shape of a Sampling Distribution

Now that you know about the mean and standard error of \bar{X} , the next step is to determine

the shape of the sampling distribution of \bar{X} ; that is, the shape of the distribution of all possible sample means (all possible values of \bar{x}) from all possible samples. You proceed differently for different conditions, which I divide into two cases: 1) the original distribution for X (the population) is normal, or has a normal distribution; and 2) the original distribution for X (the population) is *not* normal, or is unknown.

Case 1: The distribution of X is normal

If X has a normal distribution, then \bar{X} does too, no matter what the sample size n is. In the example regarding the amount of time (X) for a clerical worker to complete a task (refer to the section “Sample size and standard error”), you knew X had a normal distribution (refer to the lowest curve in Figure 11-2). If you refer to the other curves in Figure 11-2, you see the average times for samples of $n = 10$ and $n = 50$ clerical workers, respectively, also have normal distributions.



When X has a normal distribution, the sample means also always have a normal distribution, no matter what size samples you take, even if you take samples of only 2 clerical workers at a time.

The difference between the curves in Figure 11-2 is not their means or their shapes, but rather their amount of variability (how close the values in the distribution are to the mean). Results based on large samples vary less and will be more concentrated around the mean than results from small samples or results from the individuals in the population.

Case 2: The distribution of X is not normal — enter the Central Limit Theorem

If X has any distribution that is *not* normal, or if its distribution is unknown, you can't automatically say the sample mean (\bar{X}) has a normal distribution. But incredibly, you can use a normal distribution to *approximate* the distribution of \bar{X} — if the sample size is large enough. This momentous result is due to what statisticians know and love as the Central Limit Theorem.



The *Central Limit Theorem* (abbreviated *CLT*) says that if X does *not* have a normal distribution (or its distribution is unknown and hence can't be deemed to be normal), the shape of the sampling distribution of \bar{X} is *approximately* normal, as long as the sample size, n , is large enough. That is, you get an *approximate* normal distribution for the means of large samples, even if the distribution of the original values (X) is *not* normal.



Most statisticians agree that if n is at least 30, this approximation will be reasonably close in most cases, although different distribution shapes for X have different values of n that are needed. The larger the sample size (n), the closer the distribution of the sample means will be to a normal distribution.

Averaging a fair die is approximately normal

Consider the die rolling example from the earlier section “Defining a Sampling Distribution.” Notice in Figure 11-1a, the distribution of X (the population of outcomes based on millions of single rolls) is flat; the individual outcomes of each roll go from 1 to 6, and each outcome is equally likely.

Things change when you look at averages. When you roll a die a large number of times (say a sample of 50 times) and look at your outcomes, you’ll probably find about the same number of 6s as 1s (note that 6 and 1 average out to 3.5); 5s as 2s (5 and 2 also average out to 3.5); and 4s as 3s (which also average out to 3.5 — do you see a pattern here?). So if you roll a die 50 times, you have a high probability of getting an overall average that’s close to 3.5. Sometimes just by chance things won’t even out as well, but that won’t happen very often with 50 rolls.

Getting an average at the extremes with 50 rolls is a very rare event. To get an average of 1 on 50 rolls, you need all 50 rolls to be 1. How likely is that? (If it happens to you, buy a lottery ticket right away, it’s the luckiest day of your life!) The same is true for getting an average near 6.

So the chance that your average of 50 rolls is close to the middle (3.5) is highest, and the chance of it being at or close to the extremes (1 or 6) is extremely low. As for averages between 1 and 6, the probabilities get smaller as you move farther from 3.5, and the probabilities get larger as you move closer to 3.5; in particular, statisticians show that the shape of the sampling distribution of sample means in Figure 11-1b is *approximately* normal as long as the sample size is large enough. (See Chapter 9 for particulars on the shape of the normal distribution.)

Note that if you roll the die even more times, the chance of the average being close to 3.5 increases, and the sampling distribution of the sample means looks more and more like a normal distribution.

Averaging an unfair die is still approximately normal

However, sometimes the values of X don’t occur with equal probability like they do when you roll a fair die. What happens then? For example, say the die isn’t fair, and the average value for many individual rolls turns out to be 2 instead of 3.5. This means the distribution of X is skewed right (more low values like 1, 2, and 3, and fewer high values like 4, 5, and 6). But if the distribution of X (millions of individual rolls of this unfair die)

is skewed right, how does the distribution of \bar{X} (average of 50 rolls of this unfair die) end up with an approximate normal distribution?

Say that one person, Bob, is doing 50 rolls. What will the distribution of Bob's outcomes look like? Bob is more likely to get low outcomes (like 1 and 2) and less likely to get high outcomes (like 5 and 6) — the distribution of Bob's outcomes will be skewed right as well.

In fact, because Bob rolled his die a large number of times (50), the distribution of his individual outcomes has a good chance of matching the distribution of X (the outcomes from millions of rolls). However, if Bob had only rolled his die a few times (say, 6 times), he would be unlikely to even get the higher numbers like 5 and 6, and hence his distribution wouldn't look as much like the distribution of X .

If you run through the results of each of a million people like Bob who rolled this unfair die 50 times, each of their million distributions will look very similar to each other and very similar to the distribution of X . The more rolls they make each time, the closer their distributions get to the distribution of X and to each other. And here is the key: If their distributions of outcomes have a similar shape, no matter what that similar shape is, then their averages will be similar as well. Some people will get higher averages than 2 by chance, and some will get lower averages by chance, but these types of averages get less and less likely the farther you get from 2. This means you're getting an *approximate* normal distribution centered at 2.



The big deal is, it doesn't matter if you started out with a skewed distribution, or some totally wacky distribution for X . Because each of them had a large sample size (number of rolls), the distributions of each person's sample results end up looking similar, so their averages will be similar, close together, and close to a normal distribution. In fancy lingo, the distribution of \bar{X} is *approximately* normal as long as n is large enough. This is all due to the Central Limit Theorem.



In order for the CLT to work when X does *not* have a normal distribution, each person needs to roll their die enough times (that is, n must be large enough) so they have a good chance of getting all possible values of X , especially those outcomes that won't occur as often. If n is too small, some folks will not get the outcomes that have low probabilities and their means will differ from the rest by more than they should. As a result, when you put all the means together, they may not congregate around a single value. In the end, the approximate normal distribution may not show up.

Clarifying three major points about the CLT

I want to alert you to a few sources of confusion about the Central Limit Theorem before

they happen to you:

- ✓ The CLT is needed only when the distribution of X is not a normal distribution or is unknown. It is *not* needed if X started out with a normal distribution.
- ✓ The formulas for the mean and standard error of \bar{X} are *not* due to the CLT. These are just mathematical results that are always true. To see these formulas, check out the sections “The Mean of a Sampling Distribution” and “Measuring Standard Error,” earlier in this chapter.
- ✓ The n stated in the CLT refers to the size of the sample you take each time, *not* the number of samples you take. Bob rolling a die 50 times is one sample of size 50, so $n = 50$. If 10 people do it, you have 10 samples, each of size 50, and n is still 50.

Finding Probabilities for the Sample Mean

After you’ve established through the conditions addressed in case 1 or case 2 (see the previous sections) that \bar{X} has a normal or *approximately* normal distribution, you’re in luck. The normal distribution is a very friendly distribution that has a table for finding probabilities and anything else you need. For example, you can find probabilities for \bar{X} by converting the \bar{x} -value to a z -value and finding probabilities using the Z-table (provided in the appendix). (See Chapter 9 for all the details on the normal and Z-distributions.)

The general conversion formula from \bar{x} -values to z -values is:

$$z = \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}}$$

Substituting the appropriate values of the mean and standard error of \bar{X} , the conversion formula becomes:

$$z = \frac{\bar{x} - \mu_x}{\sigma_x / \sqrt{n}}$$



Don’t forget to divide by the square root of n in the denominator of z . Always divide by square root of n when the question refers to the *average* of the x -values.

Revisiting the clerical worker example from the previous section “Sample size and standard error,” suppose X is the time it takes a randomly chosen clerical worker to type and send a standard letter of recommendation. Suppose X has a normal distribution, and assume the mean is 10.5 minutes and the standard deviation 3 minutes. You take a random sample of 50 clerical workers and measure their times. What is the chance that their average time is less than 9.5 minutes?

This question translates to finding $P(\bar{X} < 9.5)$. As X has a normal distribution to start with, you know \bar{X} also has an exact (not approximate) normal distribution. Converting to z , you get:

$$z = \frac{\bar{x} - \mu_x}{\sigma_x / \sqrt{n}} = \frac{9.5 - 10.5}{3 / \sqrt{50}} = -2.36$$

So you want $P(Z < -2.36)$, which equals 0.0091 (from the Z-table in the appendix). So the chance that a random sample of 50 clerical workers average less than 9.5 minutes to complete this task is 0.91% (very small).

How do you find probabilities for \bar{X} if X is *not* normal, or unknown? As a result of the CLT, the distribution of X can be non-normal or even unknown and as long as n is large enough, you can still find *approximate* probabilities for \bar{X} using the standard normal (Z-)distribution and the process described earlier. That is, convert to a z -value and find approximate probabilities using the Z-table (in the appendix).



When you use the CLT to find a probability for \bar{X} (that is, when the distribution of X is *not* normal or is unknown), be sure to say that your answer is an *approximation*.

You also want to say the approximate answer should be close because you've got a large enough n to use the CLT. (If n is not large enough for the CLT, you can use the t -distribution in many cases — see Chapter 10.)



Beyond actual calculations, probabilities about \bar{X} can help you decide whether an assumption or a claim about a population mean is on target, based on your data. In the clerical workers example, it was assumed that the average time for all workers to type up a recommendation letter was 10.5 minutes. Your sample averaged 9.5 minutes. Because the probability that they would average less than 9.5 minutes was found to be tiny (0.0091), you either got an unusually high number of fast workers in your sample just by chance, or the assumption that the average time for all workers is 10.5 minutes was simply too high. (I'm betting on the latter.) The process of checking assumptions or challenging claims about a population is called hypothesis testing; details are in Chapter 14.

The Sampling Distribution of the Sample Proportion

The Central Limit Theorem (CLT) doesn't apply only to sample means for numerical data. You can also use it with other statistics, including sample proportions for categorical data (see Chapter 6). The *population proportion*, p , is the proportion of

individuals in the population who have a certain characteristic of interest (for example, the proportion of all Americans who are registered voters, or the proportion of all teenagers who own cellphones). The *sample proportion*, denoted \hat{p} (pronounced *p-hat*), is the proportion of individuals in the sample who have that particular characteristic; in other words, the number of individuals in the sample who have that characteristic of interest divided by the total sample size (n).

For example, if you take a sample of 100 teens and find 60 of them own cell-phones, the sample proportion of cellphone-owning teens is $\hat{p} = \frac{60}{100} = 0.60$. This section examines the sampling distribution of all possible sample proportions, \hat{p} , from samples of size n from a population.

The sampling distribution of \hat{p} has the following properties:

- ✓ Its mean, denoted by $\mu_{\hat{p}}$ (pronounced *mu sub-p-hat*), equals the population proportion, p .
- ✓ Its standard error, denoted by $\sigma_{\hat{p}}$ (say *sigma sub-p-hat*), equals:

$$\sqrt{\frac{p(1-p)}{n}}$$

(Note that because n is in the denominator, the standard error decreases as n increases.)

- ✓ Due to the CLT, its shape is *approximately* normal, provided that the sample size is large enough. Therefore you can use the normal distribution to find approximate probabilities for \hat{p} .
- ✓ The larger the sample size (n), the closer the distribution of the sample proportion is to a normal distribution.



If you are interested in the number (rather than the proportion) of individuals in your sample with the characteristic of interest, you use the binomial distribution to find probabilities for your results (see Chapter 8).



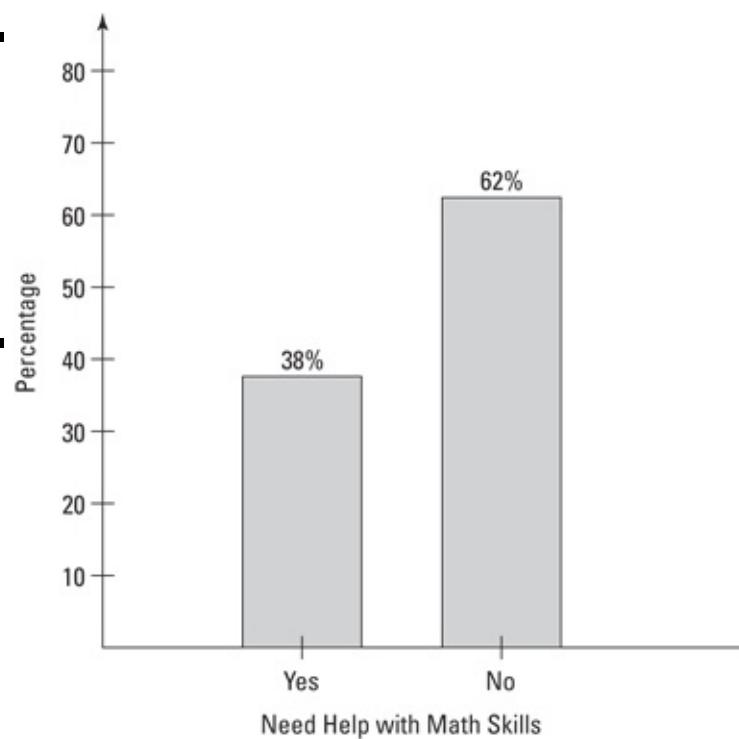
How large is large enough for the CLT to work for sample proportions? Most statisticians agree that both np and $n(1 - p)$ should be greater than or equal to 10. That is, the average number of successes (np) and the average number of failures $n(1 - p)$ needs to be at least 10.

To help illustrate the sampling distribution of the sample proportion, consider a student survey that accompanies the ACT test each year asking whether the student would like some help with math skills. Assume (through past research) that 38% of all the students taking the ACT respond yes. That means p , the population proportion, equals 0.38 in this

case. The distribution of responses (yes, no) for this population are shown in Figure 11-4 as a bar graph (see Chapter 6 for information on bar graphs).

Because 38% applies to all students taking the exam, I use p to denote the population proportion, rather than \hat{p} , which denotes sample proportions. Typically p is unknown, but I'm giving it a value here to point out how the sample proportions from samples taken from the population behave in relation to the population proportion.

Figure 11-4:
Population percentages for responses to ACT math-help question.



Now take all possible samples of $n = 1,000$ students from this population and find the proportion in each sample who said they need math help. The distribution of these sample proportions is shown in Figure 11-5. It has an *approximate* normal distribution with mean $p = 0.38$ and standard error equal to:

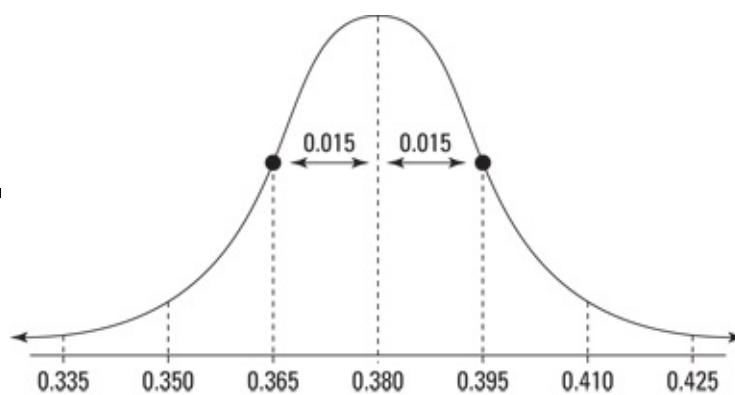
$$\sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{0.38(1-0.38)}{1,000}} = 0.015$$

(or about 1.5%).



The *approximate* normal distribution works because the two conditions for the CLT are met: 1) $np = 1,000(0.38) = 380 (\geq 10)$; and 2) $n(1 - p) = 1,000(0.62) = 620$ (also ≥ 10). And because n is so large (1,000), the approximation is excellent.

Figure 11-5:
Sampling distribution of proportion of students responding



Finding Probabilities for the Sample Proportion

You can find probabilities for \hat{p} , the sample proportion, by using the normal approximation as long as the conditions are met (see the previous section for those conditions). For the ACT test example, you assume that 0.38 or 38% of all the students taking the ACT test would like math help. Suppose you take a random sample of 100 students. What is the chance that more than 45 of them say they need math help? In terms of proportions, this is equivalent to the chance that more than $45 \div 100 = 0.45$ of them say they need help; that is, $P(\hat{p} > 0.45)$.

To answer this question, you first check the conditions: First, is np at least 10? Yes, because $100 * 0.38 = 38$. Next, is $n(1 - p)$ at least 10? Again yes, because $100 * (1 - 0.38) = 62$ checks out. So you can go ahead and use the normal approximation.

You make the conversion of the \hat{p} -value to a z -value using the following general equation:

$$z = \frac{\hat{p} - \mu_{\hat{p}}}{\sigma_{\hat{p}}} = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}}$$

When you plug in the numbers for this example, you get:

$$z = \frac{0.45 - 0.38}{\sqrt{\frac{0.38(1-0.38)}{100}}} = 1.44$$

And then you find $P(Z > 1.44) = 1 - 0.9251 = 0.0749$ using Table A-1 in the appendix. So if it's true that 0.38 percent of all students taking the exam want math help, the chance of taking a random sample of 100 students and finding more than 45 needing math help is approximately 0.0749 (by the CLT).



As noted in the previous section on sample means, you can use sample proportions to check out a claim about a population proportion. (This procedure is

a hypothesis test for a population proportion; all the details are found in Chapter 15.) In the ACT example, the probability that more than 45% of the students in a sample of 100 need math help (when you assumed 38% of the population needed math help) was found to be 0.0749. Because this probability is higher than 0.05 (the typical cutoff for blowing the whistle on a claim about a population value), you can't dispute their claim that the percentage in the population needing math help is only 38%. Our sample result is just not a rare enough event. (See Chapter 15 for more on hypothesis testing for a population proportion.)

Part IV

Guesstimating and Hypothesizing with Confidence

The 5th Wave

By Rich Tennant



"What do you mean I don't fit your desired sample population at this time?"

In this part . . .

Anytime you're given a statistic by itself, you haven't really gotten the full story. The statistic alone is missing the most important part: by how much that statistic is expected to vary. All good estimates of population parameters contain not just a statistic but also a margin of error. This combination of a statistic plus or minus a margin of error is called a confidence interval.

Now suppose you're already given a claim, assumption, or target value for the population parameter, and you want to test that claim. You do it with a hypothesis test based on sample statistics. Because sample statistics will vary, you need techniques that take this into account.

This part gives you a general, intuitive look at margin of error, confidence intervals, and hypothesis tests: their function, formulas, calculations, influential factors, and

interpretation. You also get quick references and examples for the most commonly used confidence intervals and hypothesis tests.

Chapter 12

Leaving Room for a Margin of Error

In This Chapter

- ▶ Understanding and calculating margin of error
 - ▶ Exploring the effect of sample size
 - ▶ Finding out what margin of error doesn't measure
-

Good survey and experiment researchers always include some measure of how accurate their results are so that consumers of the information can put the results into perspective. This measure is called the *margin of error (MOE)* — it's a measure of how close the sample statistic (one number that summarizes the sample) is expected to be to the population parameter being studied. (A population parameter is one number that summarizes the population. Find out more about statistics and parameters in Chapter 4.) Thankfully, many journalists are also realizing the importance of the MOE in assessing information, so reports that include the margin of error are beginning to appear in the media. But what does the margin of error really mean, and does it tell the whole story?

This chapter looks at the margin of error and what it can and can't do to help you assess the accuracy of statistical information. It also examines the issue of sample size; you may be surprised at how small a sample can be used to get a good handle on the pulse of America — or the world — if the research is done correctly.

Seeing the Importance of That Plus or Minus

Margin of error is probably not a new term to you. You've probably heard of it before, most likely in the context of survey results. For example, you may have heard someone report, "This survey had a margin of error of plus or minus three percentage points." And you may have wondered what you're supposed to do with that information and how important it really is. The truth is, the survey results themselves (with no MOE) are only a measure of how the *sample* of selected individuals felt about the issue; they don't reflect how the *entire population* may have felt, had they *all* been asked. The margin of error helps you estimate how close you are to the truth about the population based on your sample data.



Results based on a sample won't be exactly the same as what you would've found for the entire population, because when you take a sample, you don't get

information from everyone in the population. However, if the study is done right (see Chapters 16 and 17 for more about designing good studies), the results from the sample should be close to and representative of the actual values for the entire population, with a high level of confidence.



The MOE doesn't mean someone made a mistake; all it means is that you didn't get to sample everybody in the population, so you expect your sample results to vary from that population by a certain amount. In other words, you acknowledge that your results will change with subsequent samples and are only accurate to within a certain range — which can be calculated using the margin of error.

Consider one example of the type of survey conducted by some of the leading polling organizations, such as the Gallup Organization. Suppose its latest poll sampled 1,000 people from the United States, and the results show that 520 people (52%) think the president is doing a good job, compared to 48% who don't think so. Suppose Gallup reports that this survey had a margin of error of plus or minus 3%. Now, you know that the majority (more than 50%) of the people in this *sample* approve of the president, but can you say that the majority of *all Americans* approve of the president? In this case, you can't. Why not?

You need to include the margin of error (in this case, 3%) in your results. If 52% of *those sampled* approve of the president, you can expect that the percent of the *population of all Americans* who approve of the president will be 52%, plus or minus 3%. Therefore, between 49% and 55% of all Americans approve of the president. That's as close as you can get with your sample of 1,000. But notice that 49%, the lower end of this range, represents a minority, because it's less than 50%. So you really can't say that a majority of the American people support the president, based on this sample. You can only say you're confident that between 49% and 55% of all Americans support the president, which may or may not be a majority.

Think about the sample size for a moment. Isn't it interesting that a sample of only 1,000 Americans out of a population of well over 310,000,000 can lead you to be within plus or minus only 3% on your survey results? That's incredible! That means for large populations you only need to sample a tiny portion of the total to get close to the true value (assuming, as always, that you have good data). Statistics is indeed a powerful tool for finding out how people feel about issues, which is probably why so many people conduct surveys and why you're so often bothered to respond to them as well.



When you are working with categorical variables (those that record certain characteristics that don't involve measurements or counts; see Chapter 6), a quick-and-dirty way to get a rough idea of the margin of error for proportions, for any given sample size (n), is simply to find 1 divided by the square root of n . For the Gallup poll example, $n = 1,000$, and its square root is roughly 31.62, so the margin of error is

roughly 1 divided by 31.62, or about 0.03, which is equivalent to 3%. In the remainder of this chapter, you see how to get a more accurate measure of the margin of error.

Finding the Margin of Error: A General Formula

The margin of error is the amount of “plus or minus” that is attached to your sample result when you move from discussing the sample itself to discussing the whole population that it represents. Therefore, you know that the general formula for the margin of error contains a “ \pm ” in front of it. So, how do you come up with that plus or minus amount (other than taking a rough estimate, as shown above)? This section shows you how.

Measuring sample variability

Sample results vary, but by how much? According to the Central Limit Theorem (see Chapter 11), when sample sizes are large enough, the so-called sampling distribution of the sample proportions (or the sample means) follows a bell-shaped curve (or approximate normal distribution — see Chapter 9). Some of the sample proportions (or sample means) overestimate the population value and some underestimate it, but most are close to the middle.

And what’s in the middle of this sampling distribution? If you average out the results from all the possible samples you could take, the average is the actual *population proportion*, in the case of categorical data, or the actual *population average*, in the case of numerical data. Normally, you don’t know all the values of the population, so you can’t look at all of the possible sample results and average them out — but knowing something about all the other sample possibilities does help you to measure the amount by which you expect your own sample proportion (or average) to vary. (See Chapter 11 for more on sample means and proportions.)



Standard errors are the basic building blocks of the margin of error. The *standard error* of a statistic is basically equal to the standard deviation of the population divided by the square root of n (the sample size). This reflects the fact that the sample size greatly affects how much that sample statistic is going to vary from sample to sample. (See Chapter 11 for more about standard errors.)



The number of standard errors you have to add or subtract to get the MOE depends on how confident you want to be in your results (this is called your

confidence level). Typically, you want to be about 95% confident, so the basic rule is to add or subtract about 2 standard errors (1.96, to be exact) to get the MOE (you get this from the Empirical Rule; see Chapter 9). This allows you to account for about 95% of all possible results that may have occurred with repeated sampling. To be 99% confident, you add and subtract 2.58 standard errors. (This assumes a normal distribution on large n ; standard deviation known. See Chapter 11.)

You can be more precise about the number of standard errors you have to add or subtract in order to calculate the MOE for any confidence level; if the conditions are right, you can use values on the standard normal (Z -) distribution. (See Chapter 13 for details.) For any given confidence level, a corresponding value on the standard normal distribution (called a z^* -value) represents the number of standard errors to add and subtract to account for that confidence level. For 95% confidence, a more precise z^* -value is 1.96 (which is “about” 2), and for 99% confidence, the exact z^* -value is 2.58. Some of the more commonly used confidence levels (also known as percentage confidence), along with their corresponding z^* -values, are given in Table 12-1.

Table 12-1 **z^* -Values for Selected (Percentage) Confidence Levels**

Percentage Confidence	z^*-Value
80	1.28
90	1.645
95	1.96
98	2.33
99	2.58



To find a z^* -value like those in Table 12-1, add to the confidence level to make it a less-than probability and find its corresponding z -value on the Z-table. For example, a 95% confidence level means the “between” probability is 95%, so the “less-than” probability is 95% plus 2.5% (half of what’s left), or 97.5%. Look up 0.975 in the body of the Z-table and find $z^* = 1.96$ for a 95% confidence level.

Calculating margin of error for a sample proportion

When a polling question asks people to choose from a range of answers (for example, “Do you approve or disapprove the president’s performance?”), the statistic used to report the results is the proportion of people from the sample who fell into a certain group (for example, the “approve” group). This is known as the *sample proportion*. You find this number by taking the number of people in the sample that fell into the group of interest, divided by the sample size, n .

Along with the sample proportion, you need to report a margin of error. The general

formula for margin of error for the sample proportion (if certain conditions are met) is $z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$, where \hat{p} is the sample proportion, n is the sample size, and z^* is the appropriate z^* -value for your desired level of confidence (from Table 12-1). Here are the steps for calculating the margin of error for a sample proportion:

- 1. Find the sample size, n , and the sample proportion, \hat{p} .**

The sample proportion is the number in the sample with the characteristic of interest, divided by n .

- 2. Multiply the sample proportion by $(1-\hat{p})$.**

- 3. Divide the result by n .**

- 4. Take the square root of the calculated value.**

You now have the standard error, $\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$.

- 5. Multiply the result by the appropriate z^* -value for the confidence level desired.**

Refer to Table 12-1 for the appropriate z^* -value. If the confidence level is 95%, the z^* -value is 1.96.

Looking at the example involving whether Americans approve of the president, you can find the actual margin of error. First, assume you want a 95% level of confidence, so $z^* = 1.96$. The number of Americans in the sample who said they approve of the president was found to be 520. This means that the sample proportion, \hat{p} , is $520 \div 1,000 = 0.52$. (The sample size, n , was 1,000.) The margin of error for this polling question is calculated in the following way:

$$\begin{aligned} z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} &= 1.96 \sqrt{\frac{(0.52)(0.48)}{1,000}} \\ &= (1.96)(0.0158) = 0.0310 \end{aligned}$$

According to this data, you conclude with 95% confidence that 52% of all Americans approve of the president, plus or minus 3.1%.



Two conditions need to be met in order to use a z^* -value in the formula for margin of error for a sample proportion:

1. You need to be sure that $n\hat{p}$ is at least 10.
2. You need to make sure that $n(1-\hat{p})$ is at least 10.

In the preceding example of a poll on the president, $n = 1,000$, $\hat{p} = 0.52$, and $1 - \hat{p}$ is $1 - 0.52 = 0.48$. Now check the conditions: $n\hat{p} = 1,000 * 0.52 = 520$, and $n(1-\hat{p}) = 1,000 * 0.48 = 480$. Both of these numbers are at least 10, so everything is okay.

Most surveys you come across are based on hundreds or even thousands of people, so

meeting these two conditions is usually a piece of cake (unless the sample proportion is very large or very small, requiring a larger sample size to make the conditions work).



A sample proportion is the decimal version of the sample percentage. In other words, if you have a sample percentage of 5%, you must use 0.05 in the formula, not 5. To change a percentage into decimal form, simply divide by 100. After all your calculations are finished, you can change back to a percentage by multiplying your final answer by 100%.

Reporting results

Including the margin of error allows you to make conclusions beyond your sample to the population. After you calculate and interpret the margin of error, report it along with your survey results. To report the results from the president approval poll in the previous section, you say, “Based on my sample, 52% of all Americans approve of the president, plus or minus a margin of error of 3.1%. I am 95% confident in these results.”

How does a real-life polling organization report its results? Here’s an example from Gallup:

Based on the total random sample of 1,000 adults in (this) survey, we are 95% confident that the margin of error for our sampling procedure and its results is no more than ± 3.1 percentage points.

It sounds sort of like that long list of disclaimers that comes at the end of a car-leasing advertisement. But now you can understand the fine print!



Never accept the results of a survey or study without the margin of error for the study. The MOE is the only way to estimate how close the sample statistics are to the actual population parameters you’re interested in. Sample results vary, and if a different sample had been chosen, a different sample result would have been obtained; the MOE measures that amount of difference.

The next time you hear a media story about a survey or poll that was conducted, take a closer look to see if the margin of error is given; if it’s not, you should ask why. Some news outlets are getting better about reporting the margin of error for surveys, but what about other studies?

Calculating margin of error for a sample mean

When a research question asks you to estimate a parameter based on a numerical

variable (for example, “What’s the average age of teachers?”), the statistic used to help estimate the results is the average of all the responses provided by people in the sample. This is known as the *sample mean* (or average — see Chapter 5). And just like for sample proportions, you need to report a MOE for sample means.

The general formula for margin of error for the sample mean (assuming a certain condition is met) is $z^* \frac{\sigma}{\sqrt{n}}$, where σ is the population standard deviation, n is the sample size, and z^* is the appropriate z^* -value for your desired level of confidence (which you can find in Table 12-1).

Here are the steps for calculating the margin of error for a sample mean:

1. Find the population standard deviation, σ , and the sample size, n .

The population standard deviation will be given in the problem.

2. Divide the population standard deviation by the square root of the sample size.

$\frac{\sigma}{\sqrt{n}}$ gives you the standard error.

3. Multiply by the appropriate z^* -value (refer to Table 12-1).

For example, the z^* -value is 1.96 if you want to be about 95% confident.



The condition you need to meet in order to use a z^* -value in the margin of error formula for a sample mean is either: 1) The original population has a normal distribution to start with, or 2) The sample size is large enough so the normal distribution can be used (that is, the Central Limit Theorem kicks in; see Chapter 11). In general, the sample size, n , should be above about 30 for the Central Limit Theorem. Now, if it's 29, don't panic — 30 is not a magic number, it's just a general rule of thumb. (The population standard deviation must be known either way.)

Suppose you’re the manager of an ice cream shop, and you’re training new employees to be able to fill the large-size cones with the proper amount of ice cream (10 ounces each). You want to estimate the average weight of the cones they make over a one-day period, including a margin of error. Instead of weighing every single cone made, you ask each of your new employees to randomly spot check the weights of a random sample of the large cones they make and record those weights on a notepad. For $n = 50$ cones sampled, the sample mean was found to be 10.3 ounces. Suppose the population standard deviation of $\sigma = 0.6$ ounces is known.

What’s the margin of error? (Assume you want a 95% level of confidence.) It’s calculated this way:

$$z^* \frac{\sigma}{\sqrt{n}} = 1.96 \frac{0.6}{\sqrt{50}} = (1.96)(0.0849) = 0.17$$

So to report these results, you say that based on the sample of 50 cones, you estimate

that the average weight of all large cones made by the new employees over a one-day period is 10.3 ounces, with a margin of error of plus or minus 0.17 ounces. In other words, the range of likely values for the average weight of all large cones made for the day is estimated (with 95% confidence) to be between $10.30 - 0.17 = 10.13$ ounces and $10.30 + 0.17 = 10.47$ ounces. The new employees appear to be giving out too much ice cream (but I have a feeling the customers aren't offended).



Notice in the ice-cream-cone example, the units are ounces, not percentages! When working with and reporting results about data, always remember what the units are. Also, be sure that statistics are reported with their correct units of measure, and if they're not, ask what the units are.



In cases where n is too small (in general, less than 30) for the Central Limit Theorem to be used, but you still think the data came from a normal distribution, you can use a t^* -value instead of a z^* -value in your formulas. A t^* -value is one that comes from a t -distribution with $n - 1$ degrees of freedom. (Chapter 10 gives you all the in-depth details on the t -distribution.) In fact, many statisticians go ahead and use t^* -values instead of z^* -values consistently, because if the sample size is large, t^* -values and z^* -values are approximately equal anyway. In addition, for cases where you don't know the population standard deviation, σ , you can substitute it with s , the sample standard deviation; from there you use a t^* -value instead of a z^* -value in your formulas as well.

Being confident you're right

If you want to be *more* than 95% confident about your results, you need to add and subtract more than 1.96 standard errors (see Table 12-1). For example, to be 99% confident, you add and subtract 2.58 standard errors to obtain your margin of error. More confidence means a larger margin of error, though (assuming the sample size stays the same); so you have to ask yourself if it's worth it. When going from 95% to 99% confidence, the z^* -value increases by $2.58 - 1.96 = 0.62$ (see Table 12-1). Most people don't think adding and subtracting this much more of a MOE is worthwhile, just to be 4% more confident (99% versus 95%) in the results obtained.



You can never be completely certain that your sample results do reflect the population, even with the margin of error included. Even if you're 95% confident in your results, that actually means that if you repeat the sampling process over and over, 5% of the time the sample won't represent the population well, simply due to chance (not because of problems with the sampling process or anything else). In these cases, you would miss the mark. So all results need to be viewed with that in

mind.

Determining the Impact of Sample Size

The two most important ideas regarding sample size and margin of error are the following:

- ✓ Sample size and margin of error have an inverse relationship.
- ✓ After a point, increasing n beyond what you already have gives you a diminished return.

This section illustrates both concepts.

Sample size and margin of error

The relationship between margin of error and sample size is simple: As the sample size increases, the margin of error decreases. This relationship is called an inverse because the two move in opposite directions. If you think about it, it makes sense that the more information you have, the more accurate your results are going to get (in other words, the smaller your margin of error will get). (That assumes, of course, that the data were collected and handled properly.)



In the previous section, you see that the impact of a larger confidence level is a larger MOE. But if you increase the sample size, you can offset the larger MOE and bring it down to a reasonable size! Find out more about this concept in Chapter 13.

Bigger isn't always (that much) better!

In the example of the poll involving the approval rating of the president (see the earlier section “Calculating margin of error for a sample proportion”), the results of a sample of only 1,000 people from well over 310,000,000 residents in the United States could get to within about 3% of what the whole population would have said, if they had all been asked.

Using the formula for margin of error for a sample proportion, you can look at how the margin of error changes dramatically for samples of different sizes. Suppose in the presidential approval poll that n was 500 instead of 1,000. (Recall that $\hat{p} = 0.52$ for this example.) Therefore the margin of error for 95% confidence is

$$1.96 \sqrt{\frac{(0.52)(0.48)}{500}} = (1.96)(0.0223) = 0.0438,$$

which is equivalent to 4.38%. When $n = 1,000$ in the same

example, the margin of error (for 95% confidence) is $1.96\sqrt{\frac{(0.52)(0.48)}{1,000}} = (1.96)(0.0158) = 0.0310$, which is equal to 3.10%. If n is increased to 1,500, the margin of error (with the same level of confidence) becomes $1.96\sqrt{\frac{(0.52)(0.48)}{1,500}} = (1.96)(0.0129) = 0.0253$, or 2.53%. Finally, when $n = 2,000$, the margin of error is $1.96\sqrt{\frac{(0.52)(0.48)}{2,000}} = (1.96)(0.0112) = 0.0219$, or 2.19%.

Looking at these different results, you can see that larger sample sizes decrease the MOE, but after a certain point, you have a diminished return. Each time you survey one more person, the cost of your survey increases, and going from a sample size of, say, 1,500 to a sample size of 2,000 decreases your margin of error by only 0.34% (one third of one percent!) — from 0.0253 to 0.0219. The extra cost and trouble to get that small decrease in the MOE may not be worthwhile. Bigger isn't always that much better!

But what may really surprise you is that bigger can actually be worse! I explain this surprising fact in the following section.

Keeping margin of error in perspective

The margin of error is a measure of how close you expect your sample results to represent the entire population being studied. (Or at least it gives an upper limit for the amount of error you should have.) Because you're basing your conclusions about the population on your one sample, you have to account for how much those sample results could vary just due to chance.

Another view of margin of error is that it represents the maximum expected distance between the sample results and the actual population results (if you'd been able to obtain them through a census). Of course if you had the absolute truth about the population, you wouldn't be trying to do a survey, would you?

Just as important as knowing what the margin of error measures is realizing what the margin of error does *not* measure. The margin of error does not measure anything other than chance variation. That is, it doesn't measure any bias or errors that happen during the selection of the participants, the preparation or conduct of the survey, the data collection and entry process, or the analysis of the data and the drawing of the final conclusions.



A good slogan to remember when examining statistical results is “garbage in equals garbage out.” No matter how nice and scientific the margin of error may look, remember that the formula that was used to calculate it doesn't have any idea of the quality of the data that the margin of error is based on. If the sample proportion or sample mean was based on a *biased sample* (one that favored certain people over others), a bad design, bad data-collection procedures, biased questions, or

systematic errors in recording, then calculating the margin of error is pointless because it won't mean a thing.

For example, 50,000 people surveyed sounds great, but if they were all visitors to a certain Web site, the margin of error for this result is bogus because the calculation is all based on biased results! In fact, many extremely large samples are the result of biased sampling procedures. Of course, some people go ahead and report them anyway, so you have to find out what went into the formula: good information or garbage? If it turns out to be garbage, you know what to do about the margin of error. Ignore it. (For more information on errors that can take place during a survey or experiment, see Chapters 16 and 17, respectively.)

The Gallup Organization addresses the issue of what margin of error does and doesn't measure in a disclaimer that it uses to report its survey results. Gallup tells you that besides sampling error, surveys can have additional errors or bias due to question wording and some of the logistical issues involved in conducting surveys (such as missing data due to phone numbers that are no longer current).

This means that even with the best of intentions and the most meticulous attention to details and process control, stuff happens. Nothing is ever perfect. But what you need to know is that the margin of error can't measure the extent of those other types of errors. And if a highly credible polling organization like Gallup admits to possible bias, imagine what's really going on with other people's studies that aren't nearly as well designed or conducted.

Confidence Intervals: Making Your Best Guesstimate

In This Chapter

- ▶ Understanding confidence interval pieces, parts, and interpretation
 - ▶ Calculating with confidence
 - ▶ Examining factors that influence the width of a confidence interval
 - ▶ Detecting misleading results
-

Most statistics are used to estimate some characteristic about a population of interest, such as average household income, the percentage of people who buy birthday gifts online, or the average amount of ice cream consumed in the United States every year (and the resulting average weight gain — nah!). Such characteristics of a population are called *parameters*. Typically, people want to estimate (take a good guess at) the value of a parameter by taking a sample from the population and using statistics from the sample that will give them a good estimate. The question is: How do you define “good estimate”?

As long as the process is done correctly (and in the media, it often isn’t!), an estimate can often get very close to the parameter. This chapter gives you an overview of confidence intervals (the type of estimates used and recommended by statisticians); why they should be used (as opposed to just a one-number estimate); how to set up, calculate, and interpret the most commonly used confidence intervals; and how to spot misleading estimates.

Not All Estimates Are Created Equal

Read any magazine or newspaper or listen to any newscast, and you hear a number of statistics, many of which are estimates of some quantity or another. You may wonder how they came up with those statistics. In some cases, the numbers are well researched; in other cases, they’re just a shot in the dark. Here are some examples of estimates that I came across in one single issue of a leading business magazine. They come from a variety of sources:

- ✓ Even though some jobs are harder to get these days, some areas are really looking for recruits: Over the next eight years, 13,000 nurse anesthetists will be needed. Pay starts from \$80,000 to \$95,000.

- ✓ The average number of bats used by a major league baseball player per season is 90.
- ✓ The Lamborghini Murcielago can go from 0 to 60 mph in 3.7 seconds with a top speed of near 205 miles per hour.

Some of these estimates are easier to obtain than others. Here are some observations I was able to make about those estimates:

- ✓ How do you estimate how many nurse anesthetists are needed over the next eight years? You can start by looking at how many will be retiring in that time; but that won't account for growth. A prediction of the need in the next year or two would be close, but eight years into the future is much harder to do.
- ✓ The average number of bats used per major league baseball player in a season could be found by surveying the players themselves, the people who take care of their equipment, or the bat companies that supply the bats.
- ✓ Determining car speed is more difficult but could be conducted as a test with a stopwatch. And they should find the average speed of many different cars (not just one) of the same make and model, under the same driving conditions each time.



Not all statistics are created equal. To determine whether a statistic is reliable and credible, don't just take it at face value. Think about whether it makes sense and how you would go about formulating an estimate. If the statistic is really important to you, find out what process was used to come up with it. (Chapter 16 handles all the elements involving surveys, and Chapter 17 gives you the lowdown on experiments.)

Linking a Statistic to a Parameter

A *parameter* is a single number that describes a population, such as the median household income for all households in the U.S. A *statistic* is a single number that describes a sample, such as the median household income of a sample of, say, 1,200 households. You typically don't know the values of parameters of populations, so you take samples and use statistics to give your best estimates.

Suppose you want to know the percentage of vehicles in the U.S. that are pickup trucks (that's the parameter, in this case). You can't look at every single vehicle, so you take a random sample of 1,000 vehicles over a range of highways at different times of the day. You find that 7% of the vehicles in your sample are pickup trucks. Now, you don't want to say that *exactly* 7% of all vehicles on U.S. roads are pickup trucks, because you know this is only based on the 1,000 vehicles you sampled. Though you hope 7% is close to the

true percentage, you can't be sure because you based your results on a sample of vehicles, not on all the vehicles in the U.S.

So what to do? You take your sample result and add and subtract some number to indicate that you are giving a range of possible values for the population parameter, rather than just assuming the sample statistic equals the population parameter (which would not be good, although it's done in the media all the time). This number that is added to and subtracted from a statistic is called the *margin of error (MOE)*. This plus or minus (denoted by \pm) that's added to any estimate helps put the results into perspective. When you know the margin of error, you have an idea of how much the sample results could change if you took another sample.



The word *error* in *margin of error* doesn't mean a mistake was made or the quality of the data was bad. It just means the results from a sample are not exactly equal to what you would have gotten if you had used the entire population. This gap measures error due to random chance, the luck of the draw — not due to bias. (That's why minimizing bias is so important when you select your sample and collect your data; see Chapters 16 and 17.)

Getting with the Jargon

A statistic plus or minus a margin of error is called a *confidence interval*:

- ✓ The word *interval* is used because your result becomes an interval. For example, say the percentage of kids who like baseball is 40%, plus or minus 3.5%. That means the percentage of kids who like baseball is somewhere between $40\% - 3.5\% = 36.5\%$ and $40\% + 3.5\% = 43.5\%$. The lower end of the interval is your statistic minus the margin of error, and the upper end is your statistic plus the margin of error.
- ✓ With all confidence intervals, you have a certain amount of confidence in being correct (guessing the parameter) with your sample in the long run. Expressed as a percent, the amount of confidence is called the *confidence level*.

You can find formulas and examples for the most commonly used confidence intervals later in this chapter.

Following are the general steps for estimating a parameter with a confidence interval. Details on Steps 1 and 4–6 are included throughout the remainder of this chapter. Steps 2 and 3 involve sampling and data collection, which are detailed in Chapter 16 (sampling and survey data collection) and Chapter 17 (data collection from experiments).

1. Choose your confidence level and your sample size.

- 2. Select a random sample of individuals from the population.**
- 3. Collect reliable and relevant data from the individuals in the sample.**
- 4. Summarize the data into a statistic, such as a mean or proportion.**
- 5. Calculate the margin of error.**
- 6. Take the statistic plus or minus the margin of error to get your final estimate of the parameter.**

This step calculates the *confidence interval* for that parameter.

Interpreting Results with Confidence

Suppose you, a research biologist, are trying to catch a fish using a hand net, and the size of your net represents the margin of error of a confidence interval. Now say your confidence level is 95%. What does this really mean? It means that if you scoop this particular net into the water over and over again, you'll catch a fish 95% of the time. Catching a fish here means your confidence interval was correct and contains the true parameter (in this case the parameter is represented by the fish itself).

But does this mean that on any given try you have a 95% chance of catching a fish after the fact? No. Is this confusing? It certainly is. Here's the scoop (no pun intended): On a single try, say you close your eyes before you scoop your net into the water. At this point, your chances of catching a fish are 95%. But then go ahead and scoop your net through the water with your eyes still closed. *After* that's done, however, you open your eyes and see one of only two possible outcomes; you either caught a fish or you didn't; probability isn't involved anymore.

Likewise, *after* data have been collected, and the confidence interval has been calculated, you either captured the true population parameter or you didn't. So you're not saying you're 95% confident that the parameter is in your particular interval. What you are 95% confident about is the process by which random samples are selected and confidence intervals are created. (That is, 95% of the time in the long run, you'll catch a fish.)

You know that this process will result in intervals that capture the population mean 95% of the time. The other 5% of the time, the data collected in the sample just by random chance has abnormally high or low values in it and doesn't represent the population. This 5% measures errors due to random chance only and doesn't include bias.



The margin of error is meaningless if the data that went into the study were biased and/or unreliable. However, you can't tell that by looking at anyone's statistical results. My best advice is to look at how the data were collected before accepting a reported margin of error as the truth (see Chapters 16 and 17 for details on data collection issues). That means asking questions before you believe a study.

Zooming In on Width

The *width* of your confidence interval is two times the margin of error. For example, suppose the margin of error is $\pm 5\%$. A confidence interval of 7%, plus or minus 5%, goes from $7\% - 5\% = 2\%$, all the way up to $7\% + 5\% = 12\%$. So the confidence interval has a width of $12\% - 2\% = 10\%$. A simpler way to calculate this is to say that the width of the confidence interval is two times the margin of error. In this case, the width of the confidence interval is $2 * 5\% = 10\%$.



The width of a confidence interval is the distance from the lower end of the interval (statistic minus margin of error) to the upper end of the interval (statistic plus margin of error). You can always calculate the width of a confidence interval quickly by taking two times the margin of error.

The ultimate goal when making an estimate using a confidence interval is to have a narrow width, because that means you're zooming in on what the parameter is. Having to add and subtract a large margin of error only makes your result much less accurate.



So, if a small margin of error is good, is smaller even better? Not always. A narrow confidence interval is a good thing — to a point. To get an extremely narrow confidence interval, you have to conduct a much larger — and expensive — study, so a point comes where the increase in price doesn't justify the marginal difference in accuracy. Most people are pretty comfortable with a margin of error of 2% to 3% when the estimate itself is a percentage (like the percentage of women, Republicans, or smokers).

How do you go about ensuring that your confidence interval will be narrow enough? You certainly want to think about this issue before collecting your data; after the data are collected, the width of the confidence interval is set.

Three factors affect the width of a confidence interval:

- ✓ Confidence level
- ✓ Sample size
- ✓ Amount of variability in the population

Each of these three factors plays an important role in influencing the width of a confidence interval. In the following sections, you explore details of each element and how they affect width.

Choosing a Confidence Level

Every confidence interval (and every margin of error, for that matter) has a percentage associated with it that represents how confident you are that the results will capture the true population parameter, depending on the luck of the draw with your random sample. This percentage is called a *confidence level*.

A confidence level helps you account for the other possible sample results you could have gotten, when you're making an estimate of a parameter using the data from only one sample. If you want to account for 95% of the other possible results, your confidence level would be 95%.



What level of confidence is typically used by researchers? I've seen confidence levels ranging from 80% to 99%. The most common confidence level is 95%. In fact, statisticians have a saying that goes, "Why do statisticians like their jobs? Because they have to be correct only 95% of the time." (Sort of catchy, isn't it? And let's see whether forecasters beat that.)

Variability in sample results is measured in terms of number of standard errors. A *standard error* is similar to the standard deviation of a data set, only a standard error applies to sample means or sample percentages that you could have gotten if different samples were taken. (See Chapter 11 for information on standard errors.)



Standard errors are the building blocks of confidence intervals. A confidence interval is a statistic plus or minus a margin of error, and the margin of error is the number of standard errors you need to get the confidence level you want.

Every confidence level has a corresponding number of standard errors that have to be added or subtracted. This number of standard errors is called a *critical value*. In a situation where you use a Z-distribution to find the number of standard errors (as described later in this chapter), you call the critical value the *z*-value* (pronounced *z-star value*). See Table 13-1 for a list of *z**-values for some of the most common confidence levels.



As the confidence level increases, the number of standard errors increases, so the margin of error increases.

Table 13-1 z*-values for Various Confidence Levels

Confidence Level	z*-value
80%	1.28

90%	1.645 (by convention)
95%	1.96
98%	2.33
99%	2.58

If you want to be more than 95% confident about your results, you need to add and subtract more than about two standard errors. For example, to be 99% confident, you would add and subtract about two and a half standard errors to obtain your margin of error (2.58 to be exact). The higher the confidence level, the larger the z^* -value, the larger the margin of error, and the wider the confidence interval (assuming everything else stays the same). You have to pay a certain price for more confidence.

Note that I said “assuming everything else stays the same.” You can offset an increase in the margin of error by increasing the sample size. See the following section for more on this.

Factoring In the Sample Size

The relationship between margin of error and sample size is simple: As the sample size increases, the margin of error decreases, and the confidence interval gets narrower. This relationship confirms what you hope is true: The more information (data) you have, the more accurate your results are going to be. (That, of course, assumes that the information is good, credible information. See Chapter 3 for how statistics can go wrong.)



The margin of error formulas for the confidence intervals in this chapter all involve the sample size (n) in the denominator. For example, the formula for margin of error for the sample mean, $\frac{\pm z^* \sigma}{\sqrt{n}}$ (which you’ll see in great detail later in this chapter), has an n in the denominator of a fraction (this is the case for most margin of error formulas). As n increases, the denominator of this fraction increases, which makes the overall fraction get smaller. That makes the margin of error smaller and results in a narrower confidence interval.



When you need a high level of confidence, you have to increase the z^* -value and, hence, margin of error, resulting in a wider confidence interval, which isn’t good. (See the previous section.) But you can offset this wider confidence interval by increasing the sample size and bringing the margin of error back down, thus narrowing the confidence interval.

The increase in sample size allows you to still have the confidence level you want, but

also ensures that the width of your confidence interval will be small (which is what you ultimately want). You can even determine the sample size you need before you start a study: If you know the margin of error you want to get, you can set your sample size accordingly. (See the later section “Figuring Out What Sample Size You Need” for more.)



When your statistic is going to be a percentage (such as the percentage of people who prefer to wear sandals during summer), a rough way to figure margin of error for a 95% confidence interval is to take 1 divided by the square root of n (the sample size). You can try different values of n and you can see how the margin of error is affected. For example, a survey of 100 people from a large population will have a margin of error of about $\frac{1}{\sqrt{100}} = 0.10$ or plus or minus 10% (meaning the width of the confidence interval is 20%, which is pretty large).

However, if you survey 1,000 people, your margin of error decreases dramatically, to plus or minus about 3%; the width now becomes only 6%. A survey of 2,500 people results in a margin of error of plus or minus 2% (so the width is down to 4%). That’s quite a small sample size to get so accurate, when you think about how large the population is (the U.S. population, for example, is over 310 million!).

Keep in mind, however, you don’t want to go *too* high with your sample size, because a point comes where you have a diminished return. For example, moving from a sample size of 2,500 to 5,000 narrows the width of the confidence interval to about $2 * 1.4 = 2.8\%$, down from 4%. Each time you survey one more person, the cost of your survey increases, so adding another 2,500 people to the survey just to narrow the interval by little more than 1% may not be worthwhile.



The first step in any data analysis problem (and when critiquing another person’s results) is to make sure you have good data. Statistical results are only as good as the data that went into them, so real accuracy depends on the quality of the data as well as on the sample size. A large sample size that has a great deal of bias (see Chapter 16) may appear to have a narrow confidence interval — but means nothing. That’s like competing in an archery match and shooting your arrows consistently, but finding out that the whole time you’re shooting at the next person’s target; that’s how far off you are. With the field of statistics, though, you can’t accurately measure bias; you can only try to minimize it by designing good samples and studies (see Chapters 16 and 17).

Counting On Population Variability

One of the factors influencing variability in sample results is the fact that the population

itself contains variability. For example, in a population of houses in a fairly large city like Columbus, Ohio, you see a great deal of variety in not only the types of houses, but also the sizes and the prices. And the variability in prices of houses in Columbus should be more than the variability in prices of houses in a selected housing development in Columbus.

That means if you take a sample of houses from the entire city of Columbus and find the average price, the margin of error should be larger than if you take a sample from that single housing development in Columbus, even if you have the same confidence level and the same sample size.

Why? Because the houses in the entire city have more variability in price, and your sample average would change more from sample to sample than it would if you took the sample only from that single housing development, where the prices tend to be very similar because houses tend to be comparable in a single housing development. So you need to sample more houses if you're sampling from the entire city of Columbus in order to have the same amount of accuracy that you would get from that single housing development.



The standard deviation of the population is denoted σ . Notice that σ appears in the numerator of the standard error in the formula for margin of error for the sample mean: $\pm z^* \frac{\sigma}{\sqrt{n}}$.

Therefore, as the standard deviation (the numerator) increases, the standard error (the entire fraction) also increases. This results in a larger margin of error and a wider confidence interval. (Refer to Chapter 11 for more info on the standard error.)



More variability in the original population increases the margin of error, making the confidence interval wider. This increase can be offset by increasing the sample size.

Calculating a Confidence Interval for a Population Mean

When the characteristic that's being measured (such as income, IQ, price, height, quantity, or weight) is *numerical*, most people want to estimate the mean (average) value for the population. You estimate the population mean, μ , by using a sample mean, \bar{x} , plus or minus a margin of error. The result is called a *confidence interval for the population mean*, μ . Its formula depends on whether certain conditions are met. I split the conditions into two cases, illustrated in the following sections.

Case 1: Population standard deviation is known

In Case 1, the population standard deviation is known. The formula for a confidence interval (CI) for a population mean in this case is $\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$, where \bar{x} is the sample mean, σ is the population standard deviation, n is the sample size, and z^* represents the appropriate z^* -value from the standard normal distribution for your desired confidence level. (Refer to Table 13-1 for values of z^* for the given confidence levels.)



In this case, the data either have to come from a normal distribution, or if not, then n has to be large enough (at least 30 or so) for the Central Limit Theorem to kick in (see Chapter 11), allowing you to use z^* -values in the formula.

To calculate a CI for the population mean (average), under the conditions for Case 1, do the following:

1. Determine the confidence level and find the appropriate z^* -value.

Refer to Table 13-1.

2. Find the sample mean (\bar{x}) for the sample size (n).

Note: The population standard deviation is assumed to be a known value, σ .

3. Multiply z^* times σ and divide that by the square root of n .

This calculation gives you the margin of error.

4. Take \bar{x} plus or minus the margin of error to obtain the CI.

The lower end of the CI is \bar{x} minus the margin of error, whereas the upper end of the CI is \bar{x} plus the margin of error.

For example, suppose you work for the Department of Natural Resources and you want to estimate, with 95% confidence, the mean (average) length of walleye fingerlings in a fish hatchery pond.

1. Because you want a 95% confidence interval, your z^* -value is 1.96.
2. Suppose you take a random sample of 100 fingerlings and determine that the average length is 7.5 inches; assume the population standard deviation is 2.3 inches. This means $\bar{x} = 7.5$, $\sigma = 2.3$, and $n = 100$.
3. Multiply 1.96 times 2.3 divided by the square root of 100 (which is 10). The margin of error is, therefore, $\pm 1.96 * (2.3 / 10) = 1.96 * 0.23 = 0.45$ inches.
4. Your 95% confidence interval for the mean length of walleye fingerlings in this fish hatchery pond is 7.5 inches \pm 0.45 inches. (The lower end of the interval is $7.5 - 0.45 = 7.05$ inches; the upper end is $7.5 + 0.45 = 7.95$ inches.)



After you calculate a confidence interval, make sure you always interpret it in

words a non-statistician would understand. That is, talk about the results in terms of what the person in the problem is trying to find out — statisticians call this interpreting the results “in the context of the problem.” In this example you can say: “With 95% confidence, the average length of walleye fingerlings in this entire fish hatchery pond is between 7.05 and 7.95 inches, based on my sample data.” (Always be sure to include appropriate units.)

Case 2: Population standard deviation is unknown and/or n is small

In many situations, you don’t know σ , so you estimate it with the sample standard deviation, s ; and/or the sample size is small (less than 30), and you can’t be sure your data came from a normal distribution. (In the latter case, the Central Limit Theorem can’t be used; see Chapter 11.) In either situation, you can’t use a z^* -value from the standard normal (Z) distribution as your critical value anymore; you have to use a larger critical value than that, because of not knowing what σ is and/or having less data.

The formula for a confidence interval for one population mean in Case 2 is $\bar{x} \pm t_{n-1}^* \frac{s}{\sqrt{n}}$, where t_{n-1}^* is the critical t^* -value from the t -distribution with $n - 1$ degrees of freedom (where n is the sample size). The t^* -values for common confidence levels are found using the last row of the t -table (in the appendix). Chapter 10 gives you the full details on the t -distribution and how to use the t -table.



The t -distribution has a similar shape to the Z-distribution except it’s flatter and more spread out. For small values of n and a specific confidence level, the critical values on the t -distribution are larger than on the Z-distribution, so when you use the critical values from the t -distribution, the margin of error for your confidence interval will be wider. As the values of n get larger, the t^* -values are closer to z^* -values. (Chapter 10 gives you the full details on the t -distribution and its relationships to the Z-distribution.)

In the fish hatchery example from Case 1, suppose your sample size was 10 instead of 100, and everything else was the same. The t^* -value in this case comes from a t -distribution with $10 - 1 = 9$ degrees of freedom. This t^* -value is found by looking at the t -table (in the appendix). Look in the last row where the confidence levels are located, and find the confidence level of 95%; this marks the column you need. Then find the row corresponding to $df = 9$. Intersect the row and column, and you find $t^* = 2.262$. This is the t^* -value for a 95% confidence interval for the mean with a sample size of 10. (Notice this is larger than the z^* -value of 1.96 found in Table 13-1.) Calculating the confidence

interval, you get $7.5 \pm 2.262 \frac{2.3}{\sqrt{10}} = 7.50 \pm 1.645$, or 5.86 to 9.15 inches. (Chapter 10 gives you the full details on the t -distribution and how to use the t -table.)

Notice this confidence interval is wider than the one found when $n = 100$. In addition to having a larger critical value (t^* versus z^*), the sample size is much smaller, which increases the margin of error, because n is in its denominator.



In a case where you need to use s because you don't know σ , the confidence interval will be wider as well. It is also often the case that σ is unknown and the sample size is small, in which case the confidence interval is also wider.

Figuring Out What Sample Size You Need

The margin of error of a confidence interval is affected by size (see the earlier section "Factoring In the Sample Size"); as size increases, margin of error decreases. Looking at this the other way around, if you want a smaller margin of error (and doesn't everyone?), you need a larger sample size. Suppose you are getting ready to do your own survey to estimate a population mean; wouldn't it be nice to see ahead of time what sample size you need to get the margin of error you want? Thinking ahead will save you money and time and it will give you results you can live with in terms of the margin of error — you won't have any surprises later.



The formula for the sample size required to get a desired margin of error (MOE)

when you are doing a confidence interval for μ is $n \geq \left(\frac{z^* \sigma}{MOE} \right)^2$; always round up the sample size no matter what decimal value you get. (For example, if your calculations give you 126.2 people, you can't just have 0.2 of a person — you need the whole person, so include him by rounding up to 127.)

In this formula, MOE is the number representing the margin of error you want, and z^* is the z^* -value corresponding to your desired confidence level (from Table 13-1; most people use 1.96 for a 95% confidence interval). If the population standard deviation, σ , is unknown, you can put in a worst-case scenario guess for it or run a pilot study (a small trial study) ahead of time, find the standard deviation of the sample data (s), and use that number. This can be risky if the sample size is very small because it's less likely to reflect the whole population; try to get the largest trial study that you can, and/or make a conservative estimate for σ .



Often a small trial study is worth the time and effort. Not only will you get an estimate of σ to help you determine a good sample size, but you may also learn about possible problems in your data collection.



I only include one formula for calculating sample size in this chapter: the one that pertains to a confidence interval for a population mean. (You can, however, use the quick and dirty formula in the earlier section “Factoring in the Sample Size” for handling proportions.)

Here’s an example where you need to calculate n to estimate a population mean. Suppose you want to estimate the average number of songs college students store on their portable devices. You want the margin of error to be *no more than* plus or minus 20 songs. You want a 95% confidence interval. How many students should you sample?

Because you want a 95% CI, z^* is 1.96 (found in Table 13-1); you know your desired MOE is 20. Now you need a number for the population standard deviation, σ . This number is not known, so you do a pilot study of 35 students and find the standard deviation (s) for the sample is 148 songs — use this number as a substitute for σ . Using the sample size

formula, you calculate the sample size you need is $n \geq \left(\frac{1.96(148)}{20} \right)^2 = (14.504)^2 = 210.37$, which you round *up* to 211 students (you always round up when calculating n). So you need to take a random sample of *at least* 211 college students in order to have a margin of error in the number of stored songs of *no more than* 20. That’s why you see a greater-than-or-equal-to sign in the formula here.



You always round up to the nearest integer when calculating sample size, no matter what the decimal value of your result is (for example, 0.37). That’s because you want the margin of error to be *no more than* what you stated. If you round down when the decimal value is under .50 (as you normally do in other math calculations), your MOE will be a little larger than you wanted.



If you are wondering where this formula for sample size came from, it’s actually created with just a little math gymnastics. Take the margin of error formula (which contains n), fill in the remaining variables in the formula with numbers you glean from the problem, set it equal to the desired MOE, and solve for n .

Determining the Confidence Interval for One Population Proportion

When a characteristic being measured is categorical — for example, opinion on an issue (support, oppose, or are neutral), gender, political party, or type of behavior (do/don’t wear a seatbelt while driving) — most people want to estimate the proportion (or percentage) of people in the population that fall into a certain category of interest. For

example, consider the percentage of people in favor of a four-day work week, the percentage of Republicans who voted in the last election, or the proportion of drivers who don't wear seat belts. In each of these cases, the object is to estimate a population proportion, p , using a sample proportion, \hat{p} , plus or minus a margin of error. The result is called a *confidence interval for the population proportion*, p .

The formula for a CI for a population proportion is $\hat{p} \pm z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$, where \hat{p} is the sample proportion, n is the sample size, and z^* is the appropriate value from the standard normal distribution for your desired confidence level. Refer to Table 13-1 for values of z^* for certain confidence levels.

To calculate a CI for the population proportion:

1. Determine the confidence level and find the appropriate z^* -value.

Refer to Table 13-1 for z^* -values.

2. Find the sample proportion, \hat{p} , by dividing the number of people in the sample having the characteristic of interest by the sample size (n).

Note: This result should be a decimal value between 0 and 1.

3. Multiply $\hat{p}(1-\hat{p})$ and then divide that amount by n .

4. Take the square root of the result from Step 3.

5. Multiply your answer by z^* .

This step gives you the margin of error.

6. Take \hat{p} plus or minus the margin of error to obtain the CI; the lower end of the CI is \hat{p} minus the margin of error, and the upper end of the CI is \hat{p} plus the margin of error.



The formula shown in the preceding example for a CI for p is used under the condition that the sample size is large enough for the Central Limit Theorem to kick in and allow us to use a z^* -value (see Chapter 11), which happens in cases when you are estimating proportions based on large scale surveys (see Chapter 9). For small sample sizes, confidence intervals for the proportion are typically beyond the scope of an intro statistics course.

For example, suppose you want to estimate the percentage of the time you're expected to get a red light at a certain intersection.

1. Because you want a 95% confidence interval, your z^* -value is 1.96.

2. You take a random sample of 100 different trips through this intersection and find that you hit a red light 53 times, so $\hat{p} = 53/100 = 0.53$.

3. Find $\hat{p}(1-\hat{p}) = 0.53 * (1 - 0.53) = 0.2491 / 100 = 0.002491$.

4. Take the square root to get 0.0499.

The margin of error is, therefore, plus or minus $1.96 * (0.0499) = 0.0978$, or 9.78%.

5. Your 95% confidence interval for the percentage of times you will ever hit a red light at that particular intersection is 0.53 (or 53%), plus or minus 0.0978 (rounded to 0.10 or 10%). (The lower end of the interval is $0.53 - 0.10 = 0.43$ or 43%; the upper end is $0.53 + 0.10 = 0.63$ or 63%).

To interpret these results within the context of the problem, you can say that with 95% confidence the percentage of the times you should expect to hit a red light at this intersection is somewhere between 43% and 63%, based on your sample.



While performing any calculations involving sample percentages, use the decimal form. After the calculations are finished, convert to percentages by multiplying by 100. To avoid round-off error, keep at least 2 decimal places throughout.

Creating a Confidence Interval for the Difference of Two Means

The goal of many surveys and studies is to compare two populations, such as men versus women, low versus high income families, and Republicans versus Democrats. When the characteristic being compared is numerical (for example, height, weight, or income), the object of interest is the amount of difference in the means (averages) for the two populations.

For example, you may want to compare the difference in average age of Republicans versus Democrats, or the difference in average incomes of men versus women. You estimate the difference between two population means, $\mu_1 - \mu_2$, by taking a sample from each population (say, sample 1 and sample 2) and using the difference of the two sample means $\bar{x}_1 - \bar{x}_2$, plus or minus a margin of error. The result is a *confidence interval for the difference of two population means*, $\mu_1 - \mu_2$. The formula for the CI is different depending on certain conditions, as seen in the following sections; I call them Case 1 and Case 2.

Case 1: Population standard deviations are known

Case 1 assumes that both of the population standard deviations are known. The formula

for a CI for the difference between two population means (averages) is $\bar{x}_1 - \bar{x}_2 \pm z^* \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$, where \bar{x}_1 and n_1 are the mean and size of the first sample, and the first population's standard deviation, σ_1 , is given (known); \bar{x}_2 and n_2 are the mean and size of the second sample, and the second population's standard deviation, σ_2 , is given (known). Here z^* is the appropriate value from the standard normal distribution for your desired confidence level. (Refer to Table 13-1 for values of z^* for certain confidence levels.)

To calculate a CI for the difference between two population means, do the following:

1. Determine the confidence level and find the appropriate z^* -value.

Refer to Table 13-1.

2. Identify \bar{x}_1 , n_1 , and σ_1 ; find \bar{x}_2 , n_2 , and σ_2 .

3. Find the difference, $(\bar{x}_1 - \bar{x}_2)$, between the sample means.

4. Square σ_1 and divide it by n_1 ; square σ_2 and divide it by n_2 . Add the results together and take the square root.

5. Multiply your answer from Step 4 by z^* .

This answer is the margin of error.

6. Take $\bar{x}_1 - \bar{x}_2$ plus or minus the margin of error to obtain the CI.

The lower end of the CI is $\bar{x}_1 - \bar{x}_2$ minus the margin of error, whereas the upper end of the CI is $\bar{x}_1 - \bar{x}_2$ plus the margin of error.

Suppose you want to estimate with 95% confidence the difference between the mean (average) length of the cobs of two varieties of sweet corn (allowing them to grow the same number of days under the same conditions). Call the two varieties Corn-e-stats and Stats-o-sweet. Assume by prior research that the population standard deviations for Corn-e-stats and Stats-o-sweet are 0.35 inches and 0.45 inches, respectively.

1. Because you want a 95% confidence interval, your z^* is 1.96.

2. Suppose your random sample of 100 cobs of the Corn-e-stats variety averages 8.5 inches, and your random sample of 110 cobs of Stats-o-sweet averages 7.5 inches. So the information you have is: $\bar{x}_1 = 8.5$, $\sigma_1 = 0.35$, $n_1 = 100$, $\bar{x}_2 = 7.5$, $\sigma_2 = 0.45$, and $n_2 = 110$.

3. The difference between the sample means, $\bar{x}_1 - \bar{x}_2$, from Step 3, is $8.5 - 7.5 = +1$ inch. This means the average for Corn-e-stats minus the average for Stats-o-sweet is positive, making Corn-e-stats the larger of the two varieties, in terms of this sample. Is that difference enough to generalize to the entire population, though? That's what this confidence interval is going to help you decide.

4. Square σ_1 (0.35) to get 0.1225; divide by 100 to get 0.0012. Square σ_2 (0.45) and divide by 110 to get $0.2025 \div 110 = 0.0018$. The sum is $0.0012 + 0.0018 = 0.0030$; the square root is 0.0554 inches (if no rounding was done).

5. Multiply 1.96 times 0.0554 to get 0.1085 inches, the margin of error.

6. Your 95% confidence interval for the difference between the average lengths for these two varieties of sweet corn is 1 inch, plus or minus 0.1085 inches. (The lower end of the interval is $1 - 0.1085 = 0.8915$ inches; the upper end is $1 + 0.1085 = 1.1085$ inches.) Notice all the values in this interval are positive. That means Corn-e-stats is estimated to be longer than Stats-o-sweet, based on your data.

To interpret these results in the context of the problem, you can say with 95% confidence that the Corn-e-stats variety is longer, on average, than the Stats-o-sweet variety, by somewhere between 0.8915 and 1.1085 inches, based on your sample.



Notice that you could get a negative value for $\bar{x}_1 - \bar{x}_2$. For example, if you had switched the two varieties of corn, you would have gotten -1 for this difference. You would say that Stats-o-sweet averaged one inch shorter than Corn-e-stats in the sample (the same conclusion stated differently).



If you want to avoid negative values for the difference in sample means, always make the group with the larger sample mean your first group — all your differences will be positive (that's what I do).

Case 2: Population standard deviations are unknown and/or sample sizes are small

In many situations, you don't know σ_1 and σ_2 , and you estimate them with the sample standard deviations, s_1 , and s_2 ; and/or the sample sizes are small (less than 30) and you can't be sure whether your data came from a normal distribution.

A confidence interval for the difference in two population means under Case 2 is

$$(\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}},$$

where t^* is the critical value from the t -distribution with $n_1 + n_2 - 2$ degrees of freedom; n_1 and n_2 are the two sample sizes, respectively; and s_1 and s_2 are the two sample standard deviations. This t^* -value is found on the t -table (in the appendix) by intersecting the row for $df = n_1 + n_2 - 2$ with the column for the confidence level you need, as indicated by looking at the last row of the table. (See Chapter 10.) Here we assume the population standard deviations are similar; if not, modify by using the standard error and degrees of freedom. See the end of the section on comparing two means in Chapter 15.

In the corn example from Case 1, suppose the mean cob lengths of the two brands of corn, Corn-e-stats (group 1) and Stats-o-sweet (group 2), are the same as they were before: $\bar{x}_1 = 8.5$ and $\bar{x}_2 = 7.5$ inches. But this time you don't know the population standard deviations, so you use the sample standard deviations instead — suppose they turn out to be $s_1 = 0.40$ and $s_2 = 0.50$ inches, respectively. Suppose the sample sizes, n_1 and n_2 , are each only 15 in this case.

Calculating the CI, you first need to find the t^* -value on the t -distribution with $(15 + 15 - 2) = 28$ degrees of freedom. (Assume the confidence level is still 95%.) Using the t -table (in the appendix), look at the row for 28 degrees of freedom and the column representing a confidence level of 95% (see the labels on the last row of the table); intersect them and you see $t_{28}^* = 2.048$. Using the rest of the information you are given, the confidence interval for the difference in mean cob length for the two brands is

$$(8.5 - 7.5) \pm 2.048 \sqrt{\frac{(15-1)(0.4)^2 + (15-1)(0.5)^2}{15+15-2}} = 1.0 \pm 2.048(0.45) = 1.00 \pm 0.9273 \text{ inches.}$$

That means a 95% CI for the difference in the mean cob lengths of these two brands of corn in this situation is (0.0727, 1.9273) inches, with Corn-e-stats coming out on top.
(Note: This CI is wider than what was found in Case 1, as expected.)

Estimating the Difference of Two Proportions

When a characteristic, such as opinion on an issue (support/don't support), of the two groups being compared is *categorical*, people want to report on the differences between the two population proportions — for example, the difference between the proportion of women who support a four-day work week and the proportion of men who support a four-day work week. How do you do this?

You estimate the difference between two population proportions, $p_1 - p_2$, by taking a sample from each population and using the difference of the two sample proportions, $\hat{p}_1 - \hat{p}_2$, plus or minus a margin of error. The result is called a *confidence interval for the difference of two population proportions*, $p_1 - p_2$.

The formula for a CI for the difference between two population proportions is

$$(\hat{p}_1 - \hat{p}_2) \pm z^* \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}, \text{ where } \hat{p}_1 \text{ and } n_1 \text{ are the sample propor-}$$

tion and sample size of the first sample, and \hat{p}_2 and n_2 are the sample proportion and sample size of the second sample. z^* is the appropriate value from the standard normal distribution for your desired confidence level. (Refer to Table 13-1 for z^* -values.)

To calculate a CI for the difference between two population proportions, do the following:

1. Determine the confidence level and find the appropriate z^* -value.

Refer to Table 13-1.

2. Find the sample proportion \hat{p}_1 for the first sample by taking the total number from the first sample that are in the category of interest and dividing by the sample size, n_1 . Similarly, find \hat{p}_2 for the second sample.

3. Take the difference between the sample proportions, $\hat{p}_1 - \hat{p}_2$.

4. Find $\hat{p}_1(1-\hat{p}_1)$ and divide that by n_1 . Find $\hat{p}_2(1-\hat{p}_2)$ and divide that by n_2 . Add these two results together and take the square root.

5. Multiply z^* times the result from Step 4.

This step gives you the margin of error.

6. Take $\hat{p}_1 - \hat{p}_2$ plus or minus the margin of error from Step 5 to obtain the CI.

The lower end of the CI is $\hat{p}_1 - \hat{p}_2$ minus the margin of error, and the upper end of the CI is $\hat{p}_1 - \hat{p}_2$ plus the margin of error.

The formula shown here for a CI for $p_1 - p_2$ is used under the condition that both of the sample sizes are large enough for the Central Limit Theorem to kick in and allow us to use a z^* -value (see Chapter 11); this is true when you are estimating proportions using large scale surveys, for example. For small sample sizes, confidence intervals are beyond the scope of an intro statistics course.

Suppose you work for the Las Vegas Chamber of Commerce, and you want to estimate with 95% confidence the difference between the percentage of females who have ever gone to see an Elvis impersonator and the percentage of males who have ever gone to see an Elvis impersonator, in order to help determine how you should market your entertainment offerings.

1. Because you want a 95% confidence interval, your z^* -value is 1.96.
2. Suppose your random sample of 100 females includes 53 females who have seen an Elvis impersonator, so \hat{p}_1 is $53 \div 100 = 0.53$. Suppose also that your random sample of 110 males includes 37 males who have ever seen an Elvis impersonator, so \hat{p}_2 is $37 \div 110 = 0.34$.
3. The difference between these sample proportions (females – males) is $0.53 - 0.34 = 0.19$.
4. Take $0.53 * (1 - 0.53)$ and divide that by 100 to get $0.2491 \div 100 = 0.0025$. Then take $0.34 * (1 - 0.34)$ and divide that by 110 to get $0.2244 \div 110 = 0.0020$. Add these two results to get $0.0025 + 0.0020 = 0.0045$; the square root is 0.0671.
5. $1.96 * 0.0671$ gives you 0.13, or 13%, which is the margin of error.
6. Your 95% confidence interval for the difference between the percentage of females who have seen an Elvis impersonator and the percentage of males who have seen an Elvis impersonator is 0.19 or 19% (which you got in Step 3), plus or minus 13%. The lower end of the interval is $0.19 - 0.13 = 0.06$ or 6%; the upper end is $0.19 + 0.13 = 0.32$ or 32%.

To interpret these results within the context of the problem, you can say with 95% confidence that a higher percentage of females than males have seen an Elvis impersonator, and the difference in these percentages is somewhere between 6% and 32%, based on your sample.

Now I'm thinking there are some guys out there that wouldn't admit they'd ever seen an Elvis impersonator (although they've probably pretended to be one doing karaoke at some point). This may create some bias in the results. (The last time I was in Vegas, I believe I really saw Elvis; he was driving a van taxi to and from the airport. . . .)



Notice that you could get a negative value for $\hat{p}_1 - \hat{p}_2$. For example, if you had switched the males and females, you would have gotten -0.19 for this difference.

That's okay, but you can avoid negative differences in the sample proportions by having the group with the larger sample proportion serve as the first group (here, females).

Spotting Misleading Confidence Intervals

When the MOE is small, relatively speaking, you would like to say that these confidence intervals provide accurate and credible estimates of their parameters. This is not always the case, however.



Not all estimates are as accurate and reliable as the sources may want you to think. For example, a Web site survey result based on 20,000 hits may have a small MOE according to the formula, but the MOE means nothing if the survey is only given to people who happened to visit that Web site.

In other words, the sample isn't even close to being a random sample (where every sample of equal size selected from the population has an equal chance of being chosen to participate). Nevertheless, such results do get reported, along with their margins of error that make the study seem truly scientific. Beware of these bogus results! (See Chapter 12 for more on the limits of the MOE.)



Before making any decisions based on someone's estimate, do the following:

- ✓ Investigate how the statistic was created; it should be the result of a scientific process that results in reliable, unbiased, accurate data.
- ✓ Look for a margin of error. If one isn't reported, go to the original source and request it.
- ✓ Remember that if the statistic isn't reliable or contains bias, the margin of error will be meaningless.

(See Chapter 16 for evaluating survey data and see Chapter 17 for criteria for good data in experiments.)

Chapter 14

Claims, Tests, and Conclusions

In This Chapter

- ▶ Testing other people's claims
 - ▶ Using hypothesis tests to weigh evidence and make decisions
 - ▶ Recognizing that your conclusions could be wrong
-

You hear claims involving statistics all the time; the media has no shortage of them:

- ✓ Twenty-five percent of all women in the United States have varicose veins. (Wow, are some claims better left unsaid, or what?)
- ✓ Cigarette use in the U.S. continues to drop, with the percentage of all American smokers decreasing by about 2% per year over the last ten years.
- ✓ A 6-month-old baby sleeps an average of 14 to 15 hours in a 24-hour period. (Yeah, right!)
- ✓ A name-brand ready-mix pie takes only 5 minutes to make.

In today's age of information (and big money), a great deal rides on being able to back up your claims. Companies that say their products are better than the leading brand had better be able to prove it, or they could face lawsuits. Drugs that are approved by the FDA have to show strong evidence that their products actually work without producing life-threatening side effects. Manufacturers have to make sure their products are being produced according to specifications to avoid recalls, customer complaints, and loss of business.

Although many claims are backed up by solid scientific (and statistically sound) research, others are not. In this chapter, you find out how to use statistics to investigate whether a claim is actually valid and get the lowdown on the process that researchers *should* be using to validate claims that they make.



A *hypothesis test* is a statistical procedure that's designed to test a claim. Before diving into details, I want to give you the big picture of a hypothesis test by showing the main steps involved. These steps are discussed in the following sections:

- 1. Set up the null and alternative hypotheses.**
- 2. Collect good data using a well-designed study (see Chapters 16 and 17).**

3. Calculate the test statistic based on your data.
4. Find the *p*-value for your test statistic.
5. Decide whether or not to reject H_0 based on your *p*-value.
6. Understand that your conclusion may be wrong, just by chance.

Setting Up the Hypotheses

Typically in a hypothesis test, the claim being made is about a population *parameter* (one number that characterizes the entire population). Because parameters tend to be unknown quantities, everyone wants to make claims about what their values may be. For example, the claim that 25% (or 0.25) of all women have varicose veins is a claim about the proportion (that's the *parameter*) of all women (that's the *population*) who have varicose veins (that's the *variable* — having or not having varicose veins).

Researchers often challenge claims about population parameters. You may hypothesize, for example, that the actual proportion of women who have varicose veins is lower than 0.25, based on your observations. Or you may hypothesize that due to the popularity of high heeled shoes, the proportion may be higher than 0.25. Or if you're simply questioning whether the actual proportion is 0.25, your alternative hypothesis is: "No, it isn't 0.25."

Defining the null

Every hypothesis test contains a set of two opposing statements, or hypotheses, about a population parameter. The first hypothesis is called the *null hypothesis*, denoted H_0 . The null hypothesis always states that the population parameter is *equal* to the claimed value. For example, if the claim is that the average time to make a name-brand ready-mix pie is five minutes, the statistical shorthand notation for the null hypothesis in this case would be as follows: $H_0: \mu = 5$. (That is, the population mean is 5 minutes.)



All null hypotheses include an equal sign in them; there are no \leq or \geq signs in H_0 .

Not to cop out or anything, but the reason it's always equal is beyond the scope of this book; let's just say you wouldn't pay me to explain it to you.

What's the alternative?

Before actually conducting a hypothesis test, you have to put two possible hypotheses on the table — the null hypothesis is one of them. But, if the null hypothesis is rejected (that is, there was sufficient evidence against it), what's your alternative going to be?

Actually, three possibilities exist for the second (or alternative) hypothesis, denoted H_a . Here they are, along with their shorthand notations in the context of the pie example:

- ✓ The population parameter is *not equal* to the claimed value ($H_a: \mu \neq 5$).
- ✓ The population parameter is *greater than* the claimed value ($H_a: \mu > 5$).
- ✓ The population parameter is *less than* the claimed value ($H_a: \mu < 5$).

Which alternative hypothesis you choose in setting up your hypothesis test depends on what you're interested in concluding, should you have enough evidence to refute the null hypothesis (the claim).

For example, if you want to test whether a company is correct in claiming its pie takes five minutes to make and it doesn't matter whether the actual average time is more or less than that, you use the not-equal-to alternative. Your hypotheses for that test would be $H_o: \mu = 5$ versus $H_a: \mu \neq 5$.

If you only want to see whether the time turns out to be greater than what the company claims (that is, whether the company is falsely advertising its quick prep time), you use the greater-than alternative, and your two hypotheses are $H_o: \mu = 5$ versus $H_a: \mu > 5$.

Finally, say you work for the company marketing the pie, and you think the pie can be made in less than five minutes (and could be marketed by the company as such). The less-than alternative is the one you want, and your two hypotheses would be $H_o: \mu = 5$ versus $H_a: \mu < 5$.



How do you know which hypothesis to put in H_o and which one to put in H_a ?

Typically, the null hypothesis says that nothing new is happening; the previous result is the same now as it was before, or the groups have the same average (their difference is equal to zero). In general, you assume that people's claims are true until proven otherwise. So the question becomes: Can you prove otherwise? In other words, can you show sufficient evidence to reject H_o ?

Gathering Good Evidence (Data)

After you've set up the hypotheses, the next step is to collect your evidence and determine whether your evidence goes against the claim made in H_o . Remember, the claim is made about the population, but you can't test the whole population; the best you can usually do is take a sample. As with any other situation in which statistics are being collected, the quality of the data is extremely critical. (See Chapter 3 for ways to spot statistics that have gone wrong.)

Collecting good data starts with selecting a good sample. Two important issues to consider when selecting your sample are avoiding bias and being accurate. To avoid bias when selecting a sample, make it a random sample (one that's got the same chance of being selected as every other possible sample of the same size) and choose a large enough sample size so that the results will be accurate. (See Chapter 11 for more information on accuracy.)

Data is collected in many different ways, but the methods used basically boil down to two: surveys (observational studies) and experiments (controlled studies). Chapter 16 gives all the information you need to design and critique surveys, as well as information on selecting samples properly. In Chapter 17, you examine experiments: what they can do beyond an observational study, the criteria for a good experiment, and when you can conclude cause and effect.

Compiling the Evidence: The Test Statistic

After you select your sample, the appropriate number-crunching takes place. Your null hypothesis (H_0) makes a statement about the population parameter — for example, “The proportion of all women who have varicose veins is 0.25” (in other words, $H_0: p = 0.25$); or the average miles per gallon of a U.S.-built light truck is 27 ($H_0: \mu = 27$). The data you collect from the sample measures the variable of interest, and the statistics that you calculate will help you test the claim about the population parameter.

Gathering sample statistics

Say you’re testing a claim about the proportion of women with varicose veins. You need to calculate the proportion of women in your sample who have varicose veins, and that number will be your sample statistic. If you’re testing a claim about the average miles per gallon of a U.S.-built light truck, your statistic will be the average miles per gallon of the light trucks in your sample. And knowing you want to measure the variability in average miles per gallon for various trucks, you’ll want to calculate the sample standard deviation. (See Chapter 5 for all the information you need on calculating sample statistics.)

Measuring variability using standard errors

After you’ve calculated all the necessary sample statistics, you may think you’re done with the analysis part and ready to make your conclusions — but you’re not. The problem is you have no way to put your results into any kind of perspective just by looking at them in their regular units. That’s because you know that your results are based only on a sample and that sample results are going to vary. That variation needs to

be taken into account, or your conclusions could be completely wrong. (How much do sample results vary? Sample variation is measured by the standard error; see Chapter 11 for more on this.)

Suppose the claim is that the percentage of all women with varicose veins is 25%, and your sample of 100 women had 20% with varicose veins. The standard error for your sample percentage is 4% (according to formulas in Chapter 11), which means that your results are expected to vary by about twice that, or about 8%, according to the Empirical Rule (see Chapter 12). So a difference of 5%, for example, between the claim and your sample result ($25\% - 20\% = 5\%$) isn't that much, in these terms, because it represents a distance of less than 2 standard errors away from the claim.

However, suppose your sample percentage was based on a sample of 1,000 women, not 100. This decreases the amount by which you expect your results to vary, because you have more information. Again using formulas from Chapter 11, I calculate the standard error to be 0.013 or 1.3%. The margin of error (MOE) is about twice that, or 2.6% on either side. Now a difference of 5% between your sample result (20%) and the claim in H_0 (25%) is a more meaningful difference; it's way more than 2 standard errors.

Exactly how meaningful are your results? In the next section, you get more specific about measuring exactly how far apart your sample results are from the claim in terms of the number of standard errors. This leads you to a specific conclusion as to how much evidence you have against the claim in H_0 .

Understanding standard scores



The number of standard errors that a statistic lies above or below the mean is called a *standard score* (for example, a *z*-value is a type of standard score; see Chapter 9). In order to interpret your statistic, you need to convert it from original units to a standard score. When finding a standard score, you take your statistic, subtract the mean, and divide the result by the standard error.

In the case of hypothesis tests, you use the value in H_0 as the mean. (That's what you go with unless/until you have enough evidence against it.) The standardized version of your statistic is called a *test statistic*, and it's the main component of a hypothesis test. (Chapter 15 contains the formulas for the most common hypothesis tests.)

Calculating and interpreting the test statistic

The general procedure for converting a statistic to a test statistic (standard score) is as follows:

1. Take your statistic minus the claimed value (the number stated in H_0).
2. Divide by the standard error of the statistic. (Different formulas for standard error exist for different problems; see Chapter 13 for detailed formulas for standard error and Chapter 15 for formulas for various test statistics.)

Your test statistic represents the distance between your actual sample results and the claimed population value, in terms of number of standard errors. In the case of a single population mean or proportion, you know that these standardized distances should at least have an approximate standard normal distribution if your sample size is large enough (see Chapter 11). So, to interpret your test statistic in these cases, you can see where it stands on the standard normal distribution (Z-distribution).

Using the numbers from the varicose veins example in the previous section, the test statistic is found by taking the proportion in the sample with varicose veins, 0.20, subtracting the claimed proportion of all women with varicose veins, 0.25, and then dividing the result by the standard error, 0.04. These calculations give you a test statistic (standard score) of $-0.05 \div 0.04 = -1.25$. This tells you that your sample results and the population claim in H_0 are 1.25 standard errors apart; in particular, your sample results are 1.25 standard errors below the claim. Now is this enough evidence to reject the claim? The next section addresses that issue.

Weighing the Evidence and Making Decisions: p-Values

After you find your test statistic, you use it to make a decision about whether to reject H_0 . You make this decision by coming up with a number that measures the strength of this evidence (your test statistic) against the claim in H_0 . That is, how likely is it that your test statistic could have occurred while the claim was still true? This number you calculate is called the *p-value*; it's the chance that someone could have gotten results as extreme as yours while H_0 was still true. Similarly in a jury trial, the jury discusses how likely it is that all the evidence came out the way it did assuming the defendant was innocent.

This section shows all the ins and outs of *p*-values, including how to calculate them and use them to make decisions regarding H_0 .

Connecting test statistics and p-values

To test whether a claim in H_0 should be rejected (after all, it's all about H_0) you look at your test statistic taken from your sample and see whether you have enough evidence to reject the claim. If the test statistic is large (in either the positive or negative directions),

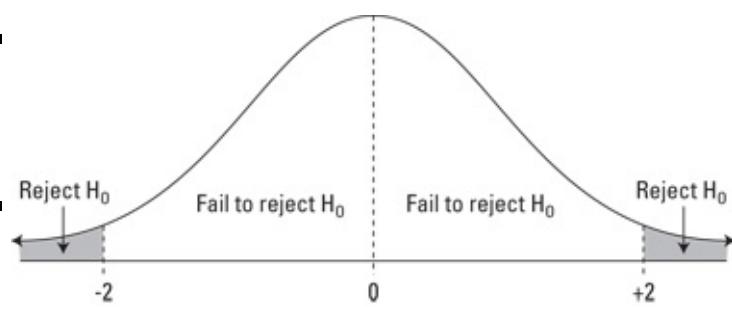
your data is far from the claim; the larger the test statistic, the more evidence you have against the claim. You determine “how far is far” by looking at where your test statistic ends up on the distribution that it came from. When testing one population mean, under certain conditions the distribution of comparison is the standard normal (*Z*-) distribution, which has a mean of 0 and a standard deviation of 1; I use it throughout this section as an example. (See Chapter 9 to find out more about the *Z*-distribution.)



If your test statistic is close to 0, or at least within that range where most of the results should fall, then you don’t have much evidence against the claim (H_0) based on your data. If your test statistic is out in the tails of the standard normal distribution (see Chapter 9 for more on tails), then your evidence against the claim (H_0) is great; this result has a very small chance of happening if the claim is true. In other words, you have sufficient evidence against the claim (H_0), and you reject H_0 .

But how far is “too far” from 0? As long as you have a normal distribution or a large enough sample size, you know that your test statistic falls somewhere on a standard normal distribution (see Chapter 11). If the null hypothesis (H_0) is true, most (about 95%) of the samples will result in test statistics that lie roughly within 2 standard errors of the claim. If H_a is the not-equal-to alternative, any test statistic outside this range will result in H_0 being rejected. See Figure 14-1 for a picture showing the locations of your test statistic and their corresponding conclusions. In the next section, you see how to quantify the amount of evidence you have against H_0 .

Figure 14-1:
Decisions for
 H_a : not-
equal-to.



Note that if the alternative hypothesis is the less-than alternative, you reject H_0 only if the test statistic falls in the left tail of the distribution (below -1.64). Similarly, if H_a is the greater-than alternative, you reject H_0 only if the test statistic falls in the right tail (above 1.64).

Defining a p-value



A *p*-value is a probability associated with your test statistic. It measures the chance of getting results at least as strong as yours if the claim (H_0) were true. In the case of

testing the population mean, the farther out your test statistic is on the tails of the standard normal (Z -) distribution, the smaller your p -value will be, the less likely your results were to have occurred, and the more evidence you have against the claim (H_0).

Calculating a p -value

To find the p -value for your test statistic:

1. Look up your test statistic on the appropriate distribution — in this case, on the standard normal (Z -) distribution (see the Z -table in the appendix).

2. Find the chance that Z is beyond (more extreme than) your test statistic:

- If H_a contains a less-than alternative, find the probability that Z is less than your test statistic (that is, look up your test statistic on the Z -table and find its corresponding probability). This is the p -value.
- If H_a contains a greater-than alternative, find the probability that Z is greater than your test statistic (look up your test statistic on the Z -table, find its corresponding probability, and subtract it from one). The result is your p -value.
- If H_a contains a non-equal-to alternative, find the probability that Z is beyond your test statistic and double it. There are two cases:

If your test statistic is negative, first find the probability that Z is less than your test statistic (look up your test statistic on the Z -table and find its corresponding probability). Then double this probability to get the p -value.

If your test statistic is positive, first find the probability that Z is greater than your test statistic (look up your test statistic on the Z -table, find its corresponding probability, and subtract it from one). Then double this result to get the p -value.



Why do you double the probabilities if your H_a contains a non-equal-to alternative? Think of the not-equal-to alternative as the combination of the greater-than alternative and the less-than alternative. If you've got a positive test statistic, its p -value only accounts for the greater-than portion of the not-equal-to alternative; double it to account for the less-than portion. (The doubling of one p -value is possible because the Z -distribution is symmetric.)

Similarly, if you've got a negative test statistic, its p -value only accounts for the less-than portion of the not-equal-to alternative; double it to also account for the greater-than portion.

When testing $H_0: p = 0.25$ versus $H_a: p < 0.25$ in the varicose veins example from the previous section, the p -value turns out to be 0.1056. This is because the test statistic (calculated in the previous section) was -1.25 , and when you look this number up on the Z -table (in the appendix) you find a probability of 0.1056 of being less than this value. If

you had been testing the two-sided alternative, $H_a: p \neq 0.25$, the p -value would be $2 * 0.1056$, or 0.2112.



If the results are likely to have occurred under the claim, then you fail to reject H_o (like a jury decides not guilty). If the results are unlikely to have occurred under the claim, then you reject H_o (like a jury decides guilty). The cutoff point between rejecting H_o and failing to reject H_o is another whole can of worms that I dissect in the next section (no pun intended).

Making Conclusions

To draw conclusions about H_o (reject or fail to reject) based on a p -value, you need to set a predetermined cutoff point where only those p -values less than or equal to the cutoff will result in rejecting H_o . This cutoff point is called the *alpha level (α)*, or *significance level* for the test. While 0.05 is a very popular cutoff value for rejecting H_o , cutoff points and resulting decisions can vary — some people use stricter cutoffs, such as 0.01, requiring more evidence before rejecting H_o , and others may have less strict cutoffs, such as 0.10, requiring less evidence.

If H_o is rejected (that is, the p -value is less than or equal to the predetermined significance level), the researcher can say she's found a statistically significant result. A result is *statistically significant* if it's too rare to have occurred by chance assuming H_o is true. If you get a statistically significant result, you have enough evidence to reject the claim, H_o , and conclude that something different or new is in effect (that is, H_a).



The significance level can be thought of as the highest possible p -value that would reject H_o and declare the results statistically significant. Following are the general rules for making a decision about H_o based on a p -value:

- ✓ If the p -value is less than or equal to your significance level, then it meets your requirements for having enough evidence against H_o ; you reject H_o .
- ✓ If the p -value is greater than your significance level, your data failed to show evidence beyond a reasonable doubt; you fail to reject H_o .

However, if you plan to make decisions about H_o by comparing the p -value to your significance level, you must decide on your significance level ahead of time. It wouldn't be fair to change your cutoff point after you've got a sneak peak at what's happening in the data.



You may be wondering whether it's okay to say "Accept H_0 " instead of "Fail to reject H_0 ." The answer is a big no. In a hypothesis test, you are *not* trying to show whether or not H_0 is true (which *accept* implies) — indeed, if you knew whether H_0 was true, you wouldn't be doing the hypothesis test in the first place. You're trying to show whether you have enough evidence to say H_0 is false, based on your data. Either you have enough evidence to say it's false (in which case you *reject* H_0) or you don't have enough evidence to say it's false (in which case you *fail to reject* H_0).

Setting boundaries for rejecting H_0 .

These guidelines help you make a decision (reject or fail to reject H_0) based on a *p*-value when your significance level is 0.05:

- ✓ If the *p*-value is less than 0.01 (very small), the results are considered highly statistically significant — *reject* H_0 .
- ✓ If the *p*-value is between 0.05 and 0.01 (but not super-close to 0.05), the results are considered statistically significant — *reject* H_0 .
- ✓ If the *p*-value is really close to 0.05 (like 0.051 or 0.049), the results should be considered marginally significant — the decision could go either way.
- ✓ If the *p*-value is greater than (but not super-close to) 0.05, the results are considered non-significant — you *fail to reject* H_0 .



When you hear a researcher say her results are found to be statistically significant, look for the *p*-value and make your own decision; the researcher's predetermined significance level may be different from yours. If the *p*-value isn't stated, ask for it.

Testing varicose veins

In the varicose veins example in the last section, the *p*-value was found to be 0.1056. This *p*-value is fairly large and indicates very weak evidence against H_0 by almost anyone's standards because it's greater than 0.05 and even slightly greater than 0.10 (considered to be a very large significance level). In this case you *fail to reject* H_0 . You didn't have enough evidence to say the proportion of women with varicose veins is less than 0.25 (your alternative hypothesis). This isn't declared to be a statistically significant result.

But say your *p*-value had been something like 0.026. A reader with a personal cutoff point of 0.05 would *reject* H_0 in this case because the *p*-value (of 0.026) is less than 0.05. His conclusion would be that the proportion of women with varicose veins isn't equal to

0.25; according to Ha in this case, you conclude it's less than 0.25, and the results are statistically significant. However, a reader whose significance level is 0.01 wouldn't have enough evidence (based on your sample) to reject H_0 because the p -value of 0.026 is greater than 0.01. These results wouldn't be statistically significant.

Finally, if the p -value turned out to be 0.049 and your significance level is 0.05, you can go by the book and say because it's less than 0.05 you reject H_0 , but you really should say your results are marginal, and let the reader decide. (Maybe they can flip a coin or something — “Heads we reject H_0 , tails, we don’t!”)

Assessing the Chance of a Wrong Decision

After you make a decision to either reject H_0 or fail to reject H_0 , the next step is living with the consequences, in terms of how people respond to your decision.

- ✓ If you conclude that a claim isn't true but it actually *is*, will that result in a lawsuit, a fine, unnecessary changes in the product, or consumer boycotts that shouldn't have happened? It's possible.
- ✓ If you can't disprove a claim that's wrong, what happens then? Will products continue to be made in the same way as they are now? Will no new law be made, no new action taken, because you showed that nothing was wrong? Missed opportunities to blow the whistle have been known to occur.



Whatever decision you make with a hypothesis test, you know there is a chance of being wrong; that's life in the statistics world. Knowing the kinds of errors that can happen and finding out how to curb the chance of them occurring are key.

Making a false alarm: Type-1 errors

Suppose a company claims that its average package delivery time is 2 days, and a consumer group tests this hypothesis, gets a p -value of 0.04, and concludes that the claim is false: They believe that the average delivery time is actually more than 2 days. This is a big deal. If the group can stand by its statistics, it has done well to inform the public about the false advertising issue. But what if the group is wrong?



Even if the group bases their study on a good design, collects good data, and makes the right analysis, it can still be wrong. Why? Because its conclusions were based on a sample of packages, not on the entire population. And as Chapter 11 tells you, sample results vary from sample to sample.

Just because the results from a sample are unusual doesn't mean they're impossible. A p -value of 0.04 means that the chance of getting your particular test statistic, even if the claim is true, is 4% (less than 5%). You reject H_0 in this case because that chance is small. But even a small chance is still a chance!

Perhaps your sample, though collected randomly, just happens to be one of those atypical samples whose result ended up far from what was expected. So, H_0 could be true, but your results lead you to a different conclusion. How often does that happen? Five percent of the time (or whatever your given cutoff probability is for rejecting H_0).



Rejecting H_0 when you shouldn't is called a *type-1 error*. I don't really like this name, because it seems so nondescript. I prefer to call a type-1 error a *false alarm*. In the case of the packages, if the consumer group made a type-1 error when it rejected the company's claim, they created a false alarm. What's the result? A very angry delivery company, I guarantee that!



To reduce the chance of false alarms, set a low cutoff probability (significance level) for rejecting H_0 . Setting it to 5% or 1% will keep the chance of a type-1 error in check.

Missing out on a detection: Type-2 errors

On the other hand, suppose the company really wasn't delivering on its claim. Who's to say that the consumer group's sample will detect it? If the actual delivery time is 2.1 days instead of 2 days, the difference would be pretty hard to detect. If the actual delivery time is 3 days, even a fairly small sample would probably show that something's up. The issue lies with those in-between values, like 2.5 days.



If H_0 is indeed false, you want to find out about it and reject H_0 . Not rejecting H_0 when you should have is called a *type-2 error*. I like to call it a *missed detection*.

Sample size is the key to being able to detect situations where H_0 is false and, thus, avoiding type-2 errors. The more information you have, the less variable your results will be (see Chapter 11) and the more ability you have to zoom in on detecting problems that exist with a claim made by H_0 .

This ability to detect when H_0 is truly false is called the *power* of a test. Power is a pretty complicated issue, but what's important for you to know is that the higher the sample size, the more powerful a test is. A powerful test has a small chance for a type-2 error.



As a preventative measure to minimize the chances of a type-2 error, statisticians recommend that you select a large sample size to ensure that any differences or departures that really exist won't be missed.

Commonly Used Hypothesis Tests: Formulas and Examples

In This Chapter

- ▶ Breaking down commonly used hypothesis tests
 - ▶ Calculating their test statistics
 - ▶ Using the results to make informed decisions
-

From product advertisements to media blitzes on recent medical breakthroughs, you often run across claims made about one or more populations. For example, “We promise to deliver our packages in two days or less” or “Two recent studies show that a high-fiber diet may reduce your risk of colon cancer by 20%.” Whenever someone makes a claim (also called a *null hypothesis*) about a population (such as all packages, or all adults) you can test the claim by doing what statisticians call a *hypothesis test*.

A hypothesis test involves setting up your *hypotheses* (a claim and its alternative), selecting a sample (or samples), collecting data, calculating the relevant statistics, and using those statistics to decide whether the claim is true.

In this chapter, I outline the formulas used for some of the most common hypothesis tests, explain the necessary calculations, and walk you through some examples.



If you need more background information on hypothesis testing (such as setting up hypotheses, understanding test statistics, p-values, significance levels, and type-1 and type-2 errors), just flip to Chapter 14. All the general concepts of hypothesis testing are developed there. This chapter focuses on their application.

Testing One Population Mean

When the variable is numerical (for example, age, income, time, and so on) and only one population or group (such as all U.S. households or all college students) is being studied, you use the hypothesis test in this section to examine or challenge a claim about the population mean. For example, a child psychologist says that the average time that working mothers spend talking to their children is 11 minutes per day, on average. (For dads, the claim is 8 minutes.) The variable — time — is numerical, and the

population is all working mothers. Using statistical notation, μ represents the average number of minutes per day that all working mothers spend talking to their children, on average.

The null hypothesis is that the population mean, μ , is equal to a certain claimed value, μ_o . The notation for the null hypothesis is $H_o: \mu = \mu_o$. So the null hypothesis in our example is $H_o: \mu = 11$ minutes, and μ_o is 11. The three possibilities for the alternative hypothesis, H_a , are $\mu \neq 11$, $\mu < 11$, or $\mu > 11$, depending on what you are trying to show. (See Chapter 14 for more on alternative hypotheses.) If you suspect that the average time working mothers spend talking with their kids is more than 11 minutes, your alternative hypothesis would be $H_a: \mu > 11$.

To test the claim, you compare the mean you got from your sample (\bar{x}) with the mean shown in H_o (μ_o). To make a proper comparison, you look at the difference between them, and divide by the standard error to take into account the fact that your sample results will vary. (See Chapter 12 for all the info you need on standard error.) This result is your *test statistic*. In the case of a hypothesis test for the population mean, the test statistic turns out (under certain conditions) to be a *z-value* (a value from the Z-distribution; see Chapter 9).

Then you can look up your test statistic on the appropriate table (in this case, you look it up on the Z-table in the appendix), and find the chance that this difference between your sample mean and the claimed population mean really could have occurred if the claim were true.

The test statistic for testing one population mean (under certain conditions) is

$$z = \frac{\bar{x} - \mu_o}{\sigma / \sqrt{n}}$$

where \bar{x} is the sample mean, σ is the population standard deviation (assume for this case that this number is known), and z is a value on the Z-distribution. To calculate the test statistic, do the following:

- 1. Calculate the sample mean, \bar{x} .**
- 2. Find $\bar{x} - \mu_o$.**
- 3. Calculate the standard error: σ / \sqrt{n} .**
- 4. Divide your result from Step 2 by the standard error found in Step 3.**



The conditions for using this test statistic are that the population standard deviation, σ , is known, and either the population has a normal distribution or the sample size is large enough to use the CLT ($n > 30$); see Chapter 11.

For our example, suppose a random sample of 100 working mothers spend an average of 11.5 minutes per day talking with their children. (Assume prior research suggests the population standard deviation is 2.3 minutes.)

1. We are given that \bar{x} is 11.5, $n = 100$, and σ is 2.3.
2. Take $11.5 - 11 = +0.5$.
3. Take 2.3 divided by the square root of 100 (which is 10) to get 0.23 for the standard error.
4. Divide +0.5 by 0.23 to get 2.17. That's your test statistic, which means your sample mean is 2.17 standard errors above the claimed population mean.



The big idea of a hypothesis test is to challenge the claim that's being made about the population (in this case, the population mean); that claim is shown in the null hypothesis, H_0 . If you have enough evidence from your sample against the claim, H_0 is rejected.

To decide whether you have enough evidence to reject H_0 , calculate the p -value by looking up your test statistic (in this case 2.17) on the standard normal (Z) distribution — see the Z -table in the appendix — and take 1 minus the probability shown. (You subtract from 1 because your H_a is a greater-than hypothesis and the table shows less-than probabilities.)

For this example you look up the test statistic (2.17) on the Z -table and find the (less-than) probability is 0.9850, so the p -value is $1 - 0.9850 = 0.015$. It's quite a bit less than your (typical) significance level 0.05, which means your sample results would be considered unusual if the claim (of 11 minutes) was true. So reject the claim ($H_0: \mu = 11$ minutes). Your results support the alternative hypothesis $H_a: \mu > 11$. According to your data, the child psychologist's claim of 11 minutes per day is too low; the actual average is greater than that.

For information on how to calculate p -values for the less-than or not-equal-to alternatives, also see Chapter 14.

Handling Small Samples and Unknown Standard Deviations: The t-Test

In two cases, you can't use the Z -distribution for a test statistic for one population mean. The first case is where the sample size is small (and by small, I mean dropping below 30 or so) the second case is when the population standard deviation, σ , is not known, and you have to estimate it using the sample standard deviation, s . In both cases, you have less reliable information on which to base your conclusions, so you have to pay a

penalty for this by using a distribution with more variability in the tails than a Z-distribution has. Enter the *t*-distribution. (See Chapter 10 for all things *t*-distribution, including its relationship with the *Z*.)

A hypothesis test for a population mean that involves the *t*-distribution is called a *t*-test. The formula for the test statistic in this case is:

$$t_{n-1} = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}, \text{ where } t_{n-1} \text{ is a value from the } t\text{-distribution with } n-1 \text{ degrees of freedom.}$$

Note it is just like the test statistic for the large sample and/or normal distribution case (see the section “Testing One Population Mean”), except σ is not known, so you substitute the sample standard deviation, s , instead, and use a *t*-value rather than a *z*-value.



Because the *t*-distribution has fatter tails than the *Z*-distribution, you get a larger *p*-value from the *t*-distribution than one that the standard normal (*Z*) distribution would have given you for the same test statistic. A bigger *p*-value means less chance of rejecting H_0 . Having less data and/or not knowing the population standard deviation should create a higher burden of proof.

Putting the *t*-test to work

Suppose a delivery company claims they deliver their packages in 2 days on average, and you suspect it’s longer than that. The hypotheses are $H_0: \mu = 2$ versus $H_a: \mu > 2$. To test this claim, you take a random sample of 10 packages and record their delivery times. You find the sample mean is $\bar{x} = 2.3$ days, and the sample standard deviation is 0.35 days. (Because the population standard deviation, σ , is unknown, you estimate it with s , the sample standard deviation.) This is a job for the *t*-test.



Because the sample size is small ($n = 10$ is much less than 30) and the population standard deviation is not known, your test statistic has a *t*-distribution. Its degrees of freedom is $10 - 1 = 9$. The formula for the test statistic (referred to as the *t*-value) is:

$$t_{10-1} = \frac{2.3 - 2.0}{0.35/\sqrt{10}} = 2.71$$

To calculate the *p*-value, you look in the row in the *t*-table (in the appendix) for $df = 9$. Your test statistic (2.71) falls between two values in the row for $df = 9$ in the *t*-table: 2.26 and 2.82 (rounding to two decimal places). To calculate the *p*-value for your test statistic, find which columns correspond to these two numbers. The number 2.26 appears in the

0.025 column and the number 2.82 appears in the 0.010 column; you now know the *p*-value for your test statistic lies between 0.025 and 0.010 (that is, $0.010 < p\text{-value} < 0.025$).

Using the *t*-table you don't know the exact number for the *p*-value, but because 0.010 and 0.025 are both less than your significance level of 0.05, you reject H_0 ; you have enough evidence in your sample to say the packages are not being delivered in 2 days, and in fact the average delivery time is more than 2 days.



The *t*-table (in the appendix) doesn't include every possible *t*-value; just find the two values closest to yours on either side, look at the columns they're in, and report your *p*-value in relation to theirs. (If your test statistic is greater than all the *t*-values in the corresponding row of the *t*-table, just use the last one; your *p*-value will be less than its probability.)



Of course you can use statistical software, if available, to calculate exact *p*-values for any test statistic; using software you get 0.012 for the exact *p*-value.

Relating *t* to *Z*

The next-to-the-last line of the *t*-table shows the corresponding values from the standard normal (*Z*) distribution for the probabilities listed on the top of each column. Now choose a column in the table and move down the column looking at the *t*-values. As the degrees of freedom of the *t*-distribution increase, the *t*-values get closer and closer to that row of the table where the *z*-values are.

This confirms a result found in Chapter 10: As the sample size (hence degrees of freedom) increases, the *t*-distribution becomes more and more like the *Z*-distribution, so the *p*-values from their hypothesis tests are virtually equal for large sample sizes. And those sample sizes don't even have to be that large to see this relationship; for $df = 30$ the *t*-values are already very similar to the *z*-values shown in the bottom of the table. These results make sense; the more data you have, the less of a penalty you have to pay. (And of course, you can use computer technology to calculate more exact *p*-values for any *t*-value you like.)

Handling negative *t*-values

For a less-than alternative hypothesis ($H_a: \bar{x} < \mu$), your test statistic would be a negative number (to the left of 0 on the *t*-distribution). In this case, you want to find the percentage below, or to the left of, your test statistic to get your *p*-value. Yet negative test statistics don't appear on the *t*-table (in the appendix).

Not to worry! The percentage to the left (below) a negative t -value is the same as the percentage to the right (above) the positive t -value, due to symmetry. So to find the p -value for your negative test statistic, look up the positive version of your test statistic on the t -table, find the corresponding right tail (greater-than) probability, and use that.

For example, suppose your test statistic is -2.7105 with 9 degrees of freedom and H_a is the less-than alternative. To find your p -value, first look up $+2.7105$ on the t -table; by the work in the previous section, you know its p -value falls between the column headings 0.025 and 0.010. Because the t -distribution is symmetric, the p -value for -2.7105 also falls somewhere between 0.025 and 0.010. Again you reject H_0 because these values are both less than or equal to 0.05.

Examining the not-equal-to alternative



To find the p -value when your alternative hypothesis (H_a) is not-equal-to, simply double the probability that you get from the t -table when you look up your test statistic. Why double it? Because the t -table shows only greater-than probabilities, which are only half the story. To find the p -value when you have a not-equal-to alternative, you must add the p -values from the less-than and greater-than alternatives. Because the t -distribution is symmetric, the less-than and greater-than probabilities are the same, so just double the one you looked up on the t -table and you'll have the p -value for the not-equal-to alternative.

For example, if your test statistic is 2.7171 and H_a is a not-equal-to alternative, look up 2.7171 on the t -table ($df = 9$ again), and you find the p -value lies between 0.025 and 0.010, as shown previously. These are the p -values for the greater-than alternative. Now double these values to include the less-than alternative and you find the p -value for your test statistic lies somewhere between $0.025 * 2 = 0.05$ and $0.010 * 2 = 0.020$.

Testing One Population Proportion

When the variable is categorical (for example, gender or support/oppose) and only one population or group is being studied (for example, all registered voters), you use the hypothesis test in this section to test a claim about the population proportion. The test looks at the proportion (p) of individuals in the population who have a certain characteristic — for example, the proportion of people who carry cellphones. The null hypothesis is $H_0: p = p_o$, where p_o is a certain claimed value of the population proportion, p . For example, if the claim is that 70% of people carry cellphones, p_o is 0.70. The alternative hypothesis is one of the following: $p > p_o$, $p < p_o$, or $p \neq p_o$. (See Chapter 14 for more on alternative hypotheses.)

The formula for the test statistic for a single proportion (under certain conditions) is:

$$z = \frac{\hat{p} - p_o}{\sqrt{\frac{p_o(1-p_o)}{n}}}$$

where \hat{p} is the proportion of individuals in the sample who have that characteristic and z is a value on the Z-distribution (see Chapter 9). To calculate the test statistic, do the following:

1. Calculate the sample proportion, \hat{p} , by taking the number of people in the sample who have the characteristic of interest (for example, the number of people in the sample carrying cellphones) and dividing that by n , the sample size.
2. Find $\hat{p} - p_o$, where p_o is the value in H_o .
3. Calculate the standard error, $\sqrt{\frac{p_o(1-p_o)}{n}}$.
4. Divide your result from Step 2 by your result from Step 3.

To interpret the test statistic, look up your test statistic on the standard normal (Z)-distribution (in the appendix) and calculate the p -value (see Chapter 14 for more on p -value calculations).



The conditions for using this test statistic are that $np_o \geq 10$ and $n(1-p_o) \geq 10$ (see Chapter 9 for details).

For example, suppose Cavifree claims that four out of five dentists recommend Cavifree toothpaste to their patients. In this case, the population is all dentists, and p is the proportion of all dentists who recommended Cavifree. The claim is that p is equal to “four out of five,” or p_o is $4 \div 5 = 0.80$. You suspect that the proportion is actually less than 0.80. Your hypotheses are $H_o: p = 0.80$ versus $H_a: p < 0.80$.

Suppose that 151 out of your sample of 200 dental patients reported receiving a recommendation for Cavifree from their dentist. To find the test statistic for these results, follow these steps:

1. Start with $\hat{p} = \frac{151}{200} = 0.755$ and $n = 200$.
2. Because $p_o = 0.80$, take $0.755 - 0.80 = -0.045$ (the numerator of the test statistic).
3. Next, the standard error equals $\sqrt{\frac{0.80(1-0.80)}{200}} = 0.028$ (the denominator of the test statistic).
4. The test statistic is $\frac{-0.045}{0.028} = -1.61$.



Because the resulting test statistic is negative, it means your sample results are –

1.61 standard errors below (less than) the claimed value for the population. How often would you expect to get results like this if H_0 were true? The chance of being at or beyond (in this case less than) -1.61 is 0.0537. (Keep the negative with the number and look up -1.61 in the Z-table in the appendix.) This result is your *p*-value because H_a is a less-than hypothesis. (See Chapter 14 for more on this.)

Because the *p*-value is greater than 0.05 (albeit not by much), you don't have quite enough evidence for rejecting H_0 . You conclude that the claim that 80% of dentists recommend Cavifree can't be rejected, according to your data. However, it's important to report the actual *p*-value too, so others can make their own decisions.



The letter *p* is used two different ways in this chapter: *p*-value and *p*. The letter *p* by itself indicates the population proportion, not the *p*-value. Don't get confused. Whenever you report a *p*-value, be sure you add *-value* so it's not confused with *p*, the population proportion.

Comparing Two (Independent) Population Averages

When the variable is numerical (for example, income, cholesterol level, or miles per gallon) and two populations or groups are being compared (for example, men versus women), you use the steps in this section to test a claim about the difference in their averages. (For example, is the difference in the population means equal to zero, indicating their means are equal?) Two independent (totally separate) random samples need to be selected, one from each population, in order to collect the data needed for this test.

The null hypothesis is that the two population means are the same; in other words, that their difference is equal to 0. The notation for the null hypothesis is $H_0: \mu_1 = \mu_2$, where μ_1 represents the mean of the first population and μ_2 represents the mean of the second population.



You can also write the null hypothesis as $H_0: \mu_1 - \mu_2 = 0$, emphasizing the idea that their difference is equal to zero if the means are the same.

The formula for the test statistic comparing two means (under certain conditions) is:

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

To calculate it, do the following:

- 1. Calculate the sample means \bar{x}_1 and \bar{x}_2 . (Assume the population standard deviations, σ_1 and σ_2 are given.) Let n_1 and n_2 represent the two sample sizes (they need not be equal).**

See Chapter 5 for these calculations.

- 2. Find the difference between the two sample means: $\bar{x}_1 - \bar{x}_2$.**



Because $\mu_1 - \mu_2$ is equal to 0 if H_0 is true, it doesn't need to be included in the numerator of the test statistic. However, if the difference they are testing is any value other than 0, you subtract that value in the numerator of the test statistic.

- 3. Calculate the standard error using the following equation:**

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

- 4. Divide your result from Step 2 by your result from Step 3.**



To interpret the test statistic, add the following two steps to the list:

- 5. Look up your test statistic on the standard normal (Z) distribution (see the Z-table in the appendix) and calculate the p -value.**

(See Chapter 14 for more on p -value calculations.)

- 6. Compare the p -value to your significance level, such as 0.05. If it's less than or equal to 0.05, reject H_0 . Otherwise, fail to reject H_0 .**

(See Chapter 14 for the details on significance levels.)



The conditions for using this test are that the two population standard deviations are known and either both populations have a normal distribution or both sample sizes are large enough for the Central Limit Theorem (see Chapter 11).

For example, suppose you want to compare the absorbency of two brands of paper towels (call the brands Stats-absorbent and Sponge-o-matic). You can make this comparison by looking at the average number of ounces each brand can absorb before being saturated. H_0 says the difference between the average absorbencies is 0 (nonexistent), and H_a says the difference is not 0. In other words, one brand is more absorbent than the other. Using statistical notation, you have $H_0 = \mu_1 - \mu_2 = 0$ versus $H_a = \mu_1 - \mu_2 \neq 0$. Here, you have no indication of which paper towel may be more absorbent, so the not-equal-to alternative is the one to use (see Chapter 14).

Suppose you select a random sample of 50 paper towels from each brand and measure the absorbency of each paper towel. Suppose the average absorbency of Stats-absorbent

(x_1) for your sample is 3 ounces, and assume the population standard deviation is 0.9 ounces. For Sponge-o-matic (x_2), the average absorbency is 3.5 ounces according to your sample; assume the population standard deviation is 1.2 ounces. Carry out this hypothesis test by following the 6 steps listed above:

- Given the above information, you know $\bar{x}_1 = 3$, $\sigma_1 = 0.9$, $\bar{x}_2 = 3.5$, $\sigma_2 = 1.2$, $n_1 = 50$, and $n_2 = 50$.
- The difference between the sample means for (Stats-absorbent – Sponge-o-matic) is $\bar{x}_1 - \bar{x}_2 = (3 - 3.5) = -0.5$ ounces. (A negative difference simply means that the second sample mean was larger than the first.)
3. The standard error is $\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = \sqrt{\frac{0.9^2}{50} + \frac{1.2^2}{50}} = \sqrt{\frac{0.81}{50} + \frac{1.44}{50}} = 0.2121$.
- Divide the difference, -0.5, by the standard error, 0.2121, which gives you -2.36. This is your test statistic.
- To find the p -value, look up -2.36 on the standard normal (Z-) distribution — see the Z-table in the appendix. The chance of being beyond, in this case to the left of, -2.36 is equal to 0.0091. Because H_a is a not-equal-to alternative, you double this percentage to get $2 * 0.0091 = 0.0182$, your p -value. (See Chapter 14 for more on the not-equal-to alternative.)
- This p -value is quite a bit less than 0.05. That means you have fairly strong evidence to reject H_o .

Your conclusion is that a statistically significant difference exists between the absorbency levels of these two brands of paper towels, based on your samples. And Sponge-o-matic comes out on top, because it has a higher average. (Stats-absorbent minus Sponge-o-matic being negative means Sponge-o-matic had the higher value.)



If one or both of your samples happen to be under 30 in size, you use the t -distribution (with degrees of freedom equal to $n_1 - 1$ or $n_2 - 1$, whichever is smaller) to look up the p -value. If the population standard deviations, σ_1 and σ_2 , are unknown, you use the sample standard deviations s_1 and s_2 instead, and you use the t -distribution with the abovementioned degrees of freedom. (See Chapter 10 for more on the t -distribution.)

Testing for an Average Difference (The Paired t-Test)

You can test for an average difference using the test in this section when the variable is numerical (for example, income, cholesterol level, or miles per gallon) and the individuals in the sample are either paired up in some way according to relevant

variables such as age or perhaps weight, or the same people are used twice (for example, using a pre-test and post-test). Paired tests are typically used for studies in which someone is testing to see whether a new treatment, technique, or method works better than an existing method, without having to worry about other factors about the subjects that may influence the results (see Chapter 17 for details).



The average difference (tested in this section) isn't the same as the difference in the averages (tested in the previous section):

- ✓ With the difference in averages, you compare the difference in the means of two separate samples to test the difference in the means of two different populations.
- ✓ With the average difference, you match up the subjects so they are thought of as coming from a single population, and the set of differences measured for each subject (for example, pre-test versus post-test) are thought of as one sample. The hypothesis test then boils down to a test for one population mean (as I explain earlier in this chapter).

For example, suppose a researcher wants to see whether teaching students to read using a computer game gives better results than teaching with a tried-and-true phonics method. She randomly selects 20 students and puts them into 10 pairs according to their reading readiness level, age, IQ, and so on. She randomly selects one student from each pair to learn to read via the computer game method (abbreviated CM), and the other learns to read using the phonics method (abbreviated PM). At the end of the study, each student takes the same reading test. The data are shown in Table 15-1.

**Table 15-1 Reading Scores for Computer Game Method
versus Phonics Method**

Student Pair	Computer Method	Phonics Method	Difference (CM – PM)
1	85	80	+5
2	80	80	0
3	95	88	+7
4	87	90	-3
5	78	72	+6
6	82	79	+3
7	57	50	+7
8	69	73	-4
9	73	78	-5
10	99	95	+4

The original data are in pairs, but you're really interested only in the difference in reading scores (computer reading score minus phonics reading score) for each pair, not

the reading scores themselves. So the *paired differences* (the differences in the pairs of scores) are your new data set. See their values in the last column of Table 15-1.

By examining the differences in the pairs of observations, you really only have a single data set, and you only have a hypothesis test for one population mean. In this case the null hypothesis is that the mean (of the paired differences) is 0, and the alternative hypothesis is that the mean (of the paired differences) is > 0 .

If the two reading methods are the same, the average of the paired differences should be 0. If the computer method is better, the average of the paired differences should be positive; the computer reading score is larger than the phonics score.



The notation for the null hypothesis is $H_0: \mu_d = 0$, where μ_d is the mean of the paired differences for the population. (The d in the subscript just reminds you that you're working with the paired differences.)

$$t_{n-1} = \frac{\bar{d} - 0}{s_d / \sqrt{n_d}}$$

The formula for the test statistic for paired differences is t_{n-1} , where \bar{d} is the average of all the paired differences found in the sample, and t_{n-1} is a value on the *t*-distribution with $n_d - 1$ degrees of freedom (see Chapter 10).



You use a *t*-distribution here because in most matched-pairs experiments the sample size is small and/or the population standard deviation σ_d is unknown, so it's estimated by s_d . (See Chapter 10 for more on the *t*-distribution.)

To calculate the test statistic for paired differences, do the following:

1. For each pair of data, take the first value in the pair minus the second value in the pair to find the paired difference.

Think of the differences as your new data set.

2. Calculate the mean, \bar{d} , and the standard deviation, s_d , of all the differences.

3. Letting nd represent the number of paired differences that you have, calculate the standard error:

$$s_d / \sqrt{n_d}$$

4. Divide \bar{d} by the standard error from Step 3.



Because μ_d is equal to 0 if H_0 is true, it doesn't really need to be included in the formula for the test statistic. As a result, you sometimes see the test statistic written like this:

$$\frac{\bar{d} - 0}{s_d / \sqrt{n_d}} = \frac{\bar{d}}{s_d / \sqrt{n_d}}$$



For the reading scores example, you can use the preceding steps to see whether the computer method is better in terms of teaching students to read.

To find the statistic, follow these steps:

1. Calculate the differences for each pair (they're shown in column 4 of Table 15-1).

Notice that the sign on each of the differences is important; it indicates which method performed better for that particular pair.

2. Calculate the mean and standard deviation of the differences from Step 1.

My calculations found the mean of the differences, $\bar{d} = 2$, and the standard deviation is $s_d = 4.64$. Note that $n_d = 10$ here.

3. The standard error is $\frac{4.64}{\sqrt{10}} = 1.47$.

(Remember that here, n_d is the number of pairs, which is 10.)

4. Take the mean of the differences (Step 2) divided by the standard error of 1.47 (Step 3) to get 1.36, the test statistic.

Is the result of Step 4 enough to say that the difference in reading scores found in this experiment applies to the whole population in general? Because the population standard deviation, σ , is unknown and you estimated it with the sample standard deviation (s), you need to use the t -distribution rather than the Z-distribution to find your p -value (see the section “Handling Small Samples and Unknown Standard Deviations: The t -Test,” earlier in this chapter). Using the t -table (in the appendix) you look up 1.36 on the t -distribution with $10 - 1 = 9$ degrees of freedom to calculate the p -value.

The p -value in this case is greater than 0.05 because 1.36 is smaller than (or to the left of) the value of 1.38 on the table, and therefore its p -value is more than 0.10 (the p -value for the column heading corresponding to 1.38).

Because the p -value is greater than 0.05, you fail to reject H_0 ; you don't have enough evidence that the mean difference in the scores between the computer method and the phonics method is significantly greater than 0. However, that doesn't necessarily mean a real difference isn't present in the population of all students. But the researcher can't say the computer game is a better reading method based on this sample of 10 students. (See Chapter 14 for information on the power of a hypothesis test and its relationship to sample size.)



In many paired experiments, the data sets are small due to costs and time associated with doing these kinds of studies. That means the t -distribution (see the t -

table in the appendix) is often used instead of the standard normal (Z-) distribution (the Z-table in the appendix) when figuring out the *p*-value.

Comparing Two Population Proportions

This test is used when the variable is categorical (for example, smoker/nonsmoker, Democrat/Republican, support/oppose an opinion, and so on) and you're interested in the proportion of individuals with a certain characteristic — for example, the proportion of smokers. In this case, two populations or groups are being compared (such as the proportion of female smokers versus male smokers).

In order to conduct this test, two independent (separate) random samples need to be selected, one from each population. The null hypothesis is that the two population proportions are the same; in other words, that their difference is equal to 0. The notation for the null hypothesis is $H_0: p_1 = p_2$, where p_1 is the proportion from the first population, and p_2 is the proportion from the second population.



Stating in H_0 that the two proportions are equal is the same as saying their difference is zero. If you start with the equation $p_1 = p_2$ and subtract p_2 from each side, you get $p_1 - p_2 = 0$. So you can write the null hypothesis either way.

The formula for the test statistic comparing two proportions (under certain conditions) is

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

where \hat{p}_1 is the proportion in the first sample with the characteristic of interest, \hat{p}_2 is the proportion in the second sample with the characteristic of interest, \hat{p} is the proportion in the combined sample (all the individuals in the first and second samples together) with the characteristic of interest, and z is a value on the Z-distribution (see Chapter 9). To calculate the test statistic, do the following:

- 1. Calculate the sample proportions \hat{p}_1 and \hat{p}_2 for each sample. Let n_1 and n_2 represent the two sample sizes (they don't need to be equal).**
- 2. Find the difference between the two sample proportions, $\hat{p}_1 - \hat{p}_2$.**
- 3. Calculate the overall sample proportion \hat{p} , the total number of individuals from both samples who have the characteristic of interest (for example, the total number of smokers, male or female, in the sample), divided by the total number of individuals from both samples ($n_1 + n_2$).**
- 4. Calculate the standard error:**

$$\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

5. Divide your result from Step 2 by your result from Step 4. This answer is your test statistic.

To interpret the test statistic, look up your test statistic on the standard normal (Z-) distribution (the Z-table in the appendix) and calculate the *p*-value, then make decisions as usual (see Chapter 14 for more on *p*-values).

Consider those drug ads that pharmaceutical companies put in magazines. The front page of an ad shows a serene picture of the sun shining, flowers blooming, people smiling — their lives changed by the drug. The company claims that its drugs can reduce allergy symptoms, help people sleep better, lower blood pressure, or fix whichever other ailment it's targeted to help. The claims may sound too good to be true, but when you turn the page to the back of the ad, you see all the fine print where the drug company justifies how it's able to make its claims. (This is typically where statistics are buried!) Somewhere in the tiny print, you'll likely find a table that shows adverse effects of the drug when compared to a *control group* (subjects who take a fake drug), for fair comparison to those who actually took the real drug (the *treatment group*; see Chapter 17 for more on this).

For example, Adderall, a drug for attention deficit hyperactivity disorder (ADHD), reported that 26 of the 374 subjects (7%) who took the drug experienced vomiting as a side effect, compared to 8 of the 210 subjects (4%) who were on a *placebo* (fake drug). Note that patients didn't know which treatment they were given. In the sample, more people on the drug experienced vomiting, but is this percentage enough to say that the entire population on the drug would experience more vomiting? You can test it to see.

In this example, you have $H_0: p_1 - p_2 = 0$ versus $H_a: p_1 - p_2 > 0$, where p_1 represents the proportion of subjects who vomited using Adderall, and p_2 represents the proportion of subjects who vomited using the placebo.



Why does H_a contain a “>” sign and not a “<” sign? H_a represents the scenario in which those taking Adderall experience more vomiting than those on the placebo — that's something the FDA (and any candidate for the drug) would want to know about. But the order of the groups is important, too. You want to set it up so the Adderall group is first, so that when you take the Adderall proportion minus the placebo proportion, you get a positive number if H_a is true. If you switch the groups, the sign would have been negative.

Now calculate the test statistic:

1. First, determine that

$$\hat{p}_1 = \frac{26}{374} = 0.070 \text{ and } \hat{p}_2 = \frac{8}{210} = 0.038$$

The sample sizes are $n_1 = 374$ and $n_2 = 210$, respectively.

2. Take the difference between these sample proportions to get $\hat{p}_1 - \hat{p}_2 = 0.070 - 0.038 = 0.032$.
3. Calculate the overall sample proportion to get $\hat{p} = \frac{26+8}{374+210} = 0.058$.
4. The standard error is $\sqrt{0.058(1-0.058)\left(\frac{1}{374} + \frac{1}{210}\right)} = 0.020$.
5. Finally, the test statistic is $0.032 \div 0.020 = 1.60$. Whew!

The p -value is the percentage chance of being at or beyond (in this case to the right of) 1.60, which is $1 - 0.9452 = 0.0548$. This p -value is just slightly greater than 0.05, so, technically, you don't have quite enough evidence to reject H_0 . That means that according to your data, vomiting is not experienced any more by those taking this drug when compared to a placebo.



A p -value that's very close to that magical but somewhat arbitrary significance level of 0.05 is what statisticians call a *marginal result*. In the preceding example, because the p -value of 0.0548 is close to the borderline between accepting and rejecting H_0 , it's generally viewed as a marginal result and should be reported as such.

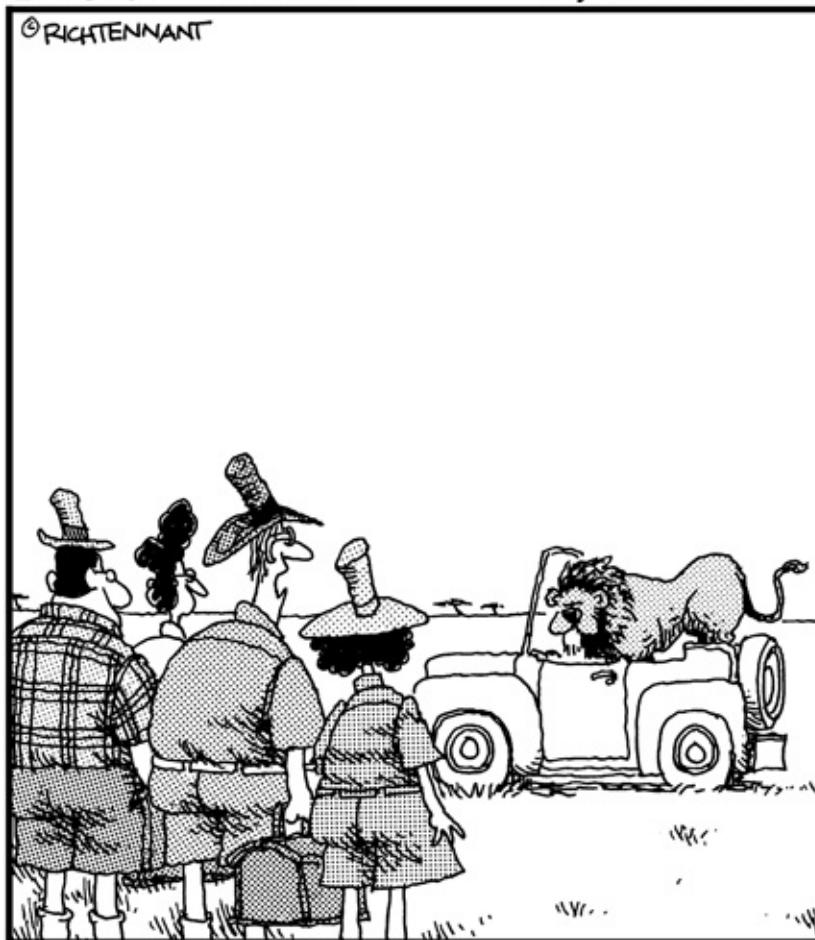
The beauty of reporting a p -value is that you can look at it and decide for yourself what you should conclude. The smaller the p -value, the more evidence you have against H_0 , but how much evidence is enough evidence? Each person is different. If you come across a report from a study in which someone found a statistically significant result, and that result is important to you, ask for the p -value so that you can make your own decision. (See Chapter 14 for more.)

Part V

Statistical Studies and the Hunt for a Meaningful Relationship

The 5th Wave

By Rich Tennant



"Okay – let's play the statistical probabilities of this situation. There are 4 of us and 1 of him. Phillip will probably start screaming, Nora will probably faint, you'll probably yell at me for leaving the truck open, and there's a good probability I'll run like a weenie if he comes toward us."

In this part . . .

Many statistics you hear and see each day are based on the results of surveys, experiments, and observational studies. Unfortunately, you can't believe everything you read or hear.

In this part, you look at what actually happens behind the scenes of these studies — how they are designed and conducted and how the data is (supposed to be) collected — so that you'll be able to spot misleading results. You also see what's needed to conduct

your own study correctly and effectively.

You also analyze data from good studies to look for relationships between two variables, where both variables are categorical (using two-way tables) or both are numerical (using correlation and regression). In addition, you see how to make proper conclusions and spot problems.

Chapter 16

Polls, Polls, and More Polls

In This Chapter

- ▶ Realizing the impact of polls and surveys
 - ▶ Going behind the scenes of polls and surveys
 - ▶ Detecting biased and inaccurate survey results
-

Surveys are all the rage amid today's information explosion. Everyone wants to know how the public feels about issues from prescription drug prices and methods of disciplining children to approval ratings of the president and ratings of reality TV shows. Polls and surveys are a big part of American life; they're a vehicle for quickly getting information about how you feel, what you think, and how you live your life, and they're a means of quickly disseminating information about important issues. Surveys highlight controversial topics, raise awareness, make political points, stress the importance of an issue, and educate or persuade the public.



Survey results can be powerful, because when many people hear that "such and such percentage of the American people do this or that," they accept these results as the truth, and then make decisions and form opinions based on that information. But in fact, many surveys *don't* provide correct, complete, or even fair or balanced information.

In this chapter, I discuss the impact of surveys and how they're used, and I take you behind the scenes of how surveys are designed and conducted so you know what to watch for when examining survey results and how to run your own surveys right. I also talk about how to interpret survey results and how to spot biased and inaccurate information, so that you can determine for yourself which results to believe and which to ignore.

Recognizing the Impact of Polls

A *survey* is an instrument that collects data through questions and answers. It is used to gather information about the opinions, behaviors, demographics, lifestyles, and other reportable characteristics of the population of interest. What's the difference between a poll and a survey? Statisticians don't make a clear distinction between the two, but I've noticed that what people call a *poll* is typically a short survey containing only a few

questions (maybe that's how researchers get more people to respond — they call it a poll rather than a survey!). But for all intents and purposes, surveys and polls are the same thing.

You come into contact with surveys and their results on a daily basis. Compared to other types of studies, such as medical experiments, some surveys can be relatively easy to conduct. They provide quick results that can often make interesting headlines in newspapers or eye-catching stories in magazines. People connect with surveys because they feel that survey results represent the opinions of people just like themselves (even though they may never have been asked to participate in a survey). And many people enjoy seeing how other people feel, what they do, where they go, and what they care about. Looking at survey results makes people feel linked with a bigger group, somehow. That's what *pollsters* (the people who conduct surveys) bank on, and that's why they spend so much time doing surveys and polls and reporting the results of this research.

Getting to the source

Who conducts surveys these days? Pretty much anyone and everyone who has a question to ask. Some of the groups that conduct polls and report the results include the following:

- ✓ News organizations
- ✓ Political parties and candidates running for office
- ✓ Professional polling organizations (such as the Gallup Organization, the Harris Poll, Zogby International, and the National Opinion Research Center [NORC])
- ✓ Representatives of magazines, TV shows, and radio programs
- ✓ Professional research organizations (like the American Medical Association, Smithsonian Institution, and Pew Research Center for the People and the Press)
- ✓ Special-interest groups (such as the National Rifle Association, Greenpeace, and American Civil Liberties Union)
- ✓ Academic researchers
- ✓ The United States government
- ✓ Joe Six-Pack (who can easily conduct his own survey on the Internet)

Ranking the worst cars of the millennium

You may be familiar with a radio show called *Car Talk* that's typically aired Saturday mornings on National Public Radio and is hosted by "Click and Clack," two brothers in Cambridge, Massachusetts, who offer wise and wacky advice to callers with strange car problems. The show's Web site regularly

offers “just for fun” surveys on a wide range of car-related topics, such as, “Who has bumper stickers on their cars, and what do they say?” One of their surveys asked the question, “What do you think was the worst car of the millennium?” Thousands upon thousands of folks responded with their votes—but, of course, these folks don’t represent all car owners. They represent only those who listen to the radio show, logged on to the Web site, and answered the survey question.

Just so you won’t be left hanging (and I know you’re dying to find out!), the results of the survey are shown in the following table. Although you may not be old enough to remember some of these vehicles, it is certainly an easy exercise to search the Internet for pictures and stories about them galore. (Remember, though, that these results represent only the opinions of *Car Talk* fans who took the time to get to the Web site and take the survey.) Notice that the percentages won’t add up to 100% because the results in the table represent only the top ten vote-getters.

Rank	Type of Car	Percentage of Votes
1	Yugo	33.7%
2	Chevy Vega	15.8%
3	Ford Pinto	12.6%
4	AMC Gremlin	8.5%
5	Chevy Chevette	7.0%
6	Renault LeCar	4.3%
7	Dodge Aspen / Plymouth Volare	4.1%
8	Cadillac Cimarron	4.0%
9	Renault Dauphine	3.6%
10	Volkswagen (VW) Bus	2.7%



Some surveys are just for fun, and others are more serious. Be sure to check the source of any serious survey in which you’re asked to participate and for which you’re given results. Groups that have a special interest in the results should either hire an independent organization to conduct (or at least to review) the survey, or they should offer copies of the survey questions to the public. Groups should also disclose in detail how the survey was designed and conducted, so that the public can make an informed decision about the credibility of the results.

Surveying what’s hot

The topics of many surveys are driven by current events, issues, and areas of interest; after all, timeliness and relevance to the public are two of the most attractive qualities of any survey. Here are just a few examples of some of the subjects being brought to the

surface by today's surveys, along with some of the results being reported:

- ✓ Does celebrity activism influence the political opinions of the American public? (Over 90% of the American public says no, according to CBS News.)
- ✓ What percentage of Americans have dated a co-worker? (A whopping 40% have, according to a career networking Web site.)
- ✓ How many patients surf the Web to find health-related information? (55% do, according to a national medical journal.)

When you read the preceding survey results, do you find yourself thinking about what the results mean to you, rather than first asking yourself whether the results are valid? Some of the preceding survey results are more valid and accurate than others, and you should think about whether to believe the results first, before accepting them without question. Nationally known polling and research organizations such as those mentioned in the previous section are credible sources, as well as journals that are *peer-reviewed* (meaning all papers published in the journal have been reviewed by others in the field and passed a certain set of standards). And the U.S. government does a good job with their data collection as well. If you are not familiar with a group conducting a survey and the results are important to you, check out the source.

Impacting lives

Whereas some surveys are just fun to look at and think about, other surveys can have a direct impact on your life or your workplace. These life-decision surveys need to be closely scrutinized before action is taken or important decisions are made. Surveys at this level can cause politicians to change or create new laws, motivate researchers to work on the latest problems, encourage manufacturers to invent new products or change business policies and practices, and influence people's behavior and ways of thinking. The following are some examples of survey results that can impact you:

- ✓ **Children's healthcare suffers:** A survey of 400 pediatricians by the Children's National Medical Center in Washington, D.C., reported that pediatricians spend, on average, only 8 to 12 minutes with each patient.
- ✓ **Teens drink more:** According to the 2009 Partnership Attitude Tracking Study, conducted by the Partnership for a Drug-Free America, the number of teens in grades 9 through 12 that use alcohol has grown by 4% (from 35% in 2008 to 39% in 2009), reversing the downward trend experienced in the ten years prior to the survey.



Always look at how researchers define the terms they're using to collect their data. In the above example, how did they define "alcohol use"? Does it count

if the teenager tried alcohol once? Does it mean they drink alcohol on a consistent basis? Results can be misleading if the range of what or who gets counted is too wide. Find out what questions were actually asked when the data was collected.

- ✓ **Crimes go unreported:** The U.S. Bureau of Justice Crime Victimization Survey concludes that only 49.4% of violent crimes were reported to police. The reasons victims gave for not reporting crimes to the police are listed in Table 16-1.

Table 16-1 Reasons Victims Didn't Report Violent Crimes

<i>Reason for Not Reporting</i>	<i>Percentage of Victims</i>
Considered it to be a personal matter	19.2%
The offender was not successful/didn't complete the crime	15.9%
Reported the crime to another official	14.7%
Didn't consider the crime to be important enough	5.5%
Didn't think police would want to be bothered	5.3%
Lack of proof	5.0%
Fear of reprisal	4.6%
Too inconvenient/time consuming to report it	3.9%
Thought police would be biased/ineffective	2.7%
Property stolen had no ID number	0.5%
Not aware that a crime occurred until later	0.4%
Other reasons	22.3%

The most frequently given reason for not reporting a violent crime to the police was that the victim considered it to be a personal matter (19.2%). Note that almost 12% of the reasons relate to perception of the reporting process itself (for example, that it would take too much time or that the police would be bothered, biased, or ineffective).



By the way, did you notice how large the “Other reasons” category is? This large, unexplained percentage indicates that the survey can be more specific and/or more research can be done regarding why crime victims don’t report crimes. Maybe the victims themselves aren’t even sure.

Behind the Scenes: The Ins and Outs of Surveys

Surveys and their results are a part of your daily experience, and you use these results to make decisions that affect your life. (Some decisions may even be life changing.) Looking at surveys with a critical eye is important. Before taking action or making decisions based on survey results, you must determine whether those results are credible, reliable, and believable. A good way to begin developing these detective skills is to go behind the scenes and see how surveys are designed, developed, implemented, and analyzed.

The survey process can be broken down into a series of ten steps:

- 1. Clarify the purpose of your survey.**
- 2. Define the target population.**
- 3. Choose the type and timing of the survey.**
- 4. Design the introduction with ethics in mind.**
- 5. Formulate the questions.**
- 6. Select the sample.**
- 7. Carry out the survey.**
- 8. Follow up, follow up, and follow up.**
- 9. Organize and analyze the data.**
- 10. Draw conclusions.**

Each step presents its own set of special issues and challenges, but each step is critical in terms of producing survey results that are fair and accurate. This sequence of steps helps you design, plan, and implement a survey, but it can also be used to critique someone else's survey, if those results are important to you.

Planning and designing a survey

The purpose of a survey is to answer questions about a target population. The *target population* is the entire group of individuals that you're interested in drawing conclusions about. In most situations, surveying the entire target population (that is, conducting a full-blown *census*) is impossible because researchers would have to spend too much time or money to do so. Usually, the best you can do is to select a sample of individuals from the target population, survey those individuals, then draw conclusions about the target population based on the data from that sample.

Sounds easy, right? Wrong. Many potential problems arise after you realize that you can't survey everyone in the entire target population. Then, after a sample is selected, many researchers aren't sure what to do to get the data they need. Unfortunately, many surveys are conducted without taking the time needed to think through these issues, resulting in errors, misleading results, and wrong conclusions. In the following sections, I give specifics for the first five steps in the survey process.

Clarifying the purpose of your survey

This sounds like it should just be common sense, but in reality, many surveys have been designed and carried out that never met their purpose, or that met only some of the objectives, but not all of them. Getting lost in the questions and forgetting what you're really trying to find out is easy to do. In stating the purpose of a survey, be as specific as possible. Think about the types of conclusions you would want to make if you were to

write a report, and let that help you determine your goals for the survey.

Lots of researchers can't see the forest for the trees. If a restaurant manager wants to determine and compare satisfaction rates for her customers, she needs to think ahead about what kinds of comparisons she wants to make and what information she wants to be able to report on. Questions that pinpoint when the customers came into the restaurant (date and time), or even what table they were at, are relevant. And if she wants to compare satisfaction rates for, say, adults versus families, she needs to ask how many people were in the party and how many were children. But if she simply asks a couple of questions on satisfaction or throws in every question she can think of, without considering in advance why she needs the information, she may end up with more questions than answers.



The more specific you can be about the purpose of the survey, the more easily you can design questions that meet your objectives, and the better off you'll be when you need to write your report.

Defining the target population

Suppose, for example, that you want to conduct a survey to determine the extent to which people send and receive personal e-mail in the workplace. You may think that the target population is e-mail users in the workplace. However, you want to determine the *extent* to which personal e-mail is used in the workplace, so you can't just ask e-mail users, or your results would be biased against those who don't use e-mail in the workplace. But should you also include those who don't even have access to a computer during their workday? (See how fast surveys can get tricky?)

The target population that probably makes the most sense here is all the people who use Internet-connected computers in the workplace. Everyone in this group at least has access to e-mail, though only some of those with access to e-mail in the workplace actually use it, and of those who use it, only some use it for personal e-mail. (And that's what you want to find out — how much they use e-mail for that purpose.)



You need to be clear in your definition of the target population. Your definition is what helps you select the proper sample, and it also guides you in your conclusions, so that you don't overgeneralize your results. If the researcher didn't clearly define the target population, this can be a sign of other problems with the survey.

Choosing the type and timing of the survey

The next step in designing your survey is to choose what type of survey is most appropriate for the situation at hand. Surveys can be done over the phone, through the

mail, with door-to-door interviews, or over the Internet. However, not every type of survey is appropriate for every situation. For example, suppose you want to determine some of the factors that relate to illiteracy in the United States. You wouldn't want to send a survey through the mail, because people who can't read won't be able to take the survey. In that case, a telephone interview is more appropriate.



Choose the type of survey that's most appropriate for the target population, in terms of getting the most truthful and informative data possible. You also have to keep in mind the budget you have to work with; door-to-door interviews are more expensive than phone surveys, for example. When examining the results of a survey, be sure to look at whether the type of survey used is most appropriate for the situation, keeping budget considerations in mind.

Next you need to decide when to conduct the survey. In life, timing is everything, and the same goes for surveys. Current events shape people's opinions all the time, and although some pollsters try to determine how people feel about those events, others take advantage of events, especially negative ones, and use them as political platforms or as fodder for headlines and controversy. For example, surveys about gun control often come up after a shooting takes place. Also take note of other events that were going on at the time of the survey; for example, people may not want to answer their phones during the Super Bowl, on election night, during the Olympics, or around holidays. Improper timing can lead to bias.

In addition to the date, the time of day is also important. If you conduct a telephone survey to get people's opinions on stress in the workplace and you call them at home between the hours of 9 a.m. and 5 p.m., you're going to have bias in your results; those are the hours when the majority of people are at work (busy being stressed out!).

Designing the introduction with ethics in mind

While this rule doesn't apply to little polls that you see on the Internet and in magazines, serious surveys need to provide information pertaining to important ethical issues. First, they should include what pollsters call a *cover letter* — an introduction that explains the purpose of the survey, what will be done with the data, whether the information the respondent supplies will be confidential or anonymous (see the sidebar "Anonymity versus confidentiality" later in this chapter), and that the person's participation is appreciated but not required. The cover letter should also provide the researcher's contact information for respondents to use if they have questions or concerns.



If the survey is done by any institution or group that is federally regulated, such as a university, research institute, or a hospital, the survey has to be approved in advance by a committee designated to review, regulate, and/or monitor the research to make sure it's ethical, scientific, and follows regulations. Such committees are

called institutional review boards (IRBs), independent ethics committees (IECs), or ethical review boards (ERBs). The survey cover letter should explain who has approved the research. If you don't see such information, ask.

Formulating the questions

After the purpose, type, timing, and ethical issues of the survey have been addressed, the next step is to formulate the questions. The way that the questions are asked can make a huge difference in the quality of the data that will be collected. One of the single most common sources of bias in surveys is the wording of the questions. Research shows that the wording of the questions can directly affect the outcome of a survey. *Leading questions*, also called *misleading questions*, are designed to favor a certain response over another. They can greatly affect how people answer the questions, and their responses may not accurately reflect how they truly feel about an issue.

For example, here are two ways that I've seen survey questions worded about a proposed school bond issue (both of which are leading questions):

Don't you agree that a tiny percentage increase in sales tax is a worthwhile investment in improving the quality of the education of our children?

Don't you think we should stop increasing the burden on the taxpayers and stop asking for yet another sales tax hike to fund the wasteful school system?

From the wording of each of these leading questions, you can easily see how the pollsters want you to respond. To make matters worse, neither question tells you exactly how much of a tax increase is being proposed, which is also misleading.



The best way to word a question is in a neutral way, giving the reader the necessary information required to make an informed decision. For example, the tax issue question is better worded this way:

The school district is proposing a 0.01% increase in sales tax to provide funds for a new high school to be built in the district. What's your opinion on the proposed sales tax? (Possible responses: strongly in favor, in favor, neutral, against, strongly against.)

If the purpose of a survey is purely to collect information rather than influence or persuade the respondent, the questions should be worded in a neutral and informative way in order to minimize bias. The best way to assess the neutrality of a question is to ask yourself whether you can tell how the person wants you to respond. If the answer is yes, that question is a leading question and can give misleading results.



If the results of a survey are important to you, ask the researcher for a copy of the

questions used on the survey so you can assess the quality of the questions. When conducting your own survey, have others check the questions to verify that the wording is neutral and informative.

Selecting the sample

After the survey has been designed, the next step is to select people to participate in the survey. Because typically you don't have time or money to conduct a census (a survey of the entire target population), you need to select a subset of the population, called a *sample*. How this sample is selected can make all the difference in terms of the accuracy and the quality of the results.

Three criteria are important in selecting a good sample, as you find out in the following sections.

A good sample represents the target population

To represent the target population, the sample must be selected from the target population, the whole target population, and nothing but the target population. Suppose you want to find out how many hours of TV Americans watch in a day, on average. Asking students in a dorm at a local university to record their TV viewing habits isn't going to cut it. Students represent only a portion of the target population.



Unfortunately, many people who conduct surveys don't take the time or spend the money to select a representative sample of people to participate in the study, and they end up with biased survey results. When presented with survey results, find out how the sample was selected before examining the results of the survey and see how well they match the target population.

A good sample is selected randomly

A *random* sample is one in which every possible sample (of the same size) has an equal chance of being selected from the target population. The easiest example to visualize here is that of a hat (or bucket) containing individual slips of paper, each with the name of a person written on it; if the slips are thoroughly mixed before each slip of paper is drawn out, the result will be a random sample of the target population (in this case, the population of people whose names are in the hat). A random sample eliminates bias in the sampling process.

Reputable polling organizations, such as the Gallup Organization, use a random digit-dialing procedure to telephone the members of their sample. Of course, this excludes people without telephones, but because most American households today have at least one telephone, the bias involved in excluding people without telephones is relatively

small.



Beware of surveys that have a large but not randomly selected sample. Internet surveys are the biggest culprit. Someone can say that 50,000 people logged on to a Web site to answer a survey, and that means the person posting this site has gotten a lot of data. But the information is biased; research shows that people who respond to surveys tend to have stronger opinions than those that don't respond. And if they didn't even select the participants randomly to start with, imagine how strong (and biased) the respondents' opinions would be. If the survey designer sampled fewer people but did so randomly, the survey results would be more accurate.

A good sample is large enough for the results to be accurate

If you have a large sample size, and if the sample is representative of the target population and is selected at random, you can count on that information being pretty accurate. *How* accurate depends on the sample size, but the bigger the sample size, the more accurate the information will be (as long as that information is good information). The accuracy of most survey questions is measured in terms of a percentage. This percentage is called the *margin of error*, and it represents how much the researcher expects the results to vary if she were to repeat the survey many times using different samples of the same size. Read more about this in Chapter 12.



A quick and dirty formula to estimate the minimum amount of accuracy of a survey involving categorical data (such as gender or political affiliation) is to take 1 divided by the square root of the sample size. For example, a survey of 1,000 (randomly selected) people is accurate to within ± 0.032 , or 3.2 percentage points. (See Chapter 12 for the exact formula for calculating the accuracy of a survey.) In cases where not everyone responded, you should replace the sample size with the number of respondents (see the “Following up, following up, and following up” section later in this chapter). Remember, these quick-and-dirty estimates of accuracy are conservative; using the precise formulas gives you accuracy rates that are often much better than these. (See Chapter 13 for details.)



With large populations (in the thousands, say) it's the size of the sample, not the size of the population, that matters. For example, if you randomly sample 1,000 individuals from a large population, your accuracy level is estimated to be within 3.2 percentage points, no matter whether you sample from a small town of 10,000 people, a state of 1,000,000 people, or all of the United States. That fact was one of the things that blew my mind about statistics when I first learned it, and it still does today — it's amazing how accurate you can get with such a comparatively small sample size.



However, with small populations, you have to apply different methods to determine accuracy and sample size. A sample of 10 out of a population of 100 takes a much larger piece out of the pie than a sample of 10 out of 10,000 does, for example. More advanced methods involving a finite population correction handle issues that come up with small populations.

Carrying out a survey

The survey has been designed, and the participants have been selected. Now you have to go about the process of carrying out the survey, which is another important step — one where lots of mistakes and biases can occur.

Collecting the data

During the survey itself, the participants can have problems understanding the questions, they may give answers that aren't among the choices (in the case of a multiple choice question), or they may decide to give answers that are inaccurate or blatantly false; the latter is called *response bias*. (As an example of response bias, think about the difficulties involved in getting people to tell the truth about whether they've cheated on their income-tax forms.)

Some of the potential problems with the data-collection process can be minimized or avoided with careful training of the personnel who carry out the survey. With proper training, any issues that arise during the survey are resolved in a consistent and clear way, and fewer errors are made in recording the data. Problems with confusing questions or incomplete choices for answers can be resolved by conducting a pilot study on a few participants prior to the actual survey and then, based on their feedback, fixing any problems with the questions.

Personnel can also be trained to create an environment in which each respondent feels safe enough to tell the truth; ensuring that privacy will be protected also helps encourage more people to respond. To minimize interviewer bias, the interviewers must follow a script that's the same for each subject.

Anonymity versus confidentiality

If you were to conduct a survey to determine the extent of personal e-mail use at work, the response rate would probably be an issue, because many people are reluctant to discuss their use of personal e-mail in the workplace, or at least to do so truthfully. You could try to encourage people to respond by letting them know that their privacy would be protected during and after the survey.

When you report the results of a survey, you generally don't tie the information collected to the names of the respondents, because doing so would violate the privacy of the respondents. You've probably

heard the terms *anonymous* and *confidential* before, but what you may not realize is that these two words are completely different in terms of privacy issues. Keeping results *confidential* means that I could tie your information to your name in my report, but I promise that I won't do that. Keeping results *anonymous* means that I have no way of tying your information to your name in my report, even if I wanted to.

If you're asked to participate in a survey, be sure you're clear about what the researchers plan to do with your responses and whether or not your name can be tied to the survey. (Good surveys always make this issue very clear for you.) Then make a decision as to whether you still want to participate.



Beware of conflicts of interest that come up with misleading surveys. For example, if you are being asked about the quality of your service by the person who gave you the service, you may not want to respond truthfully. Or, if your physical therapist gives you an “anonymous” feedback survey on your last day and tells you to give it to her when you’re done, the survey may have issues of bias.

Following up, following up, and following up

Anyone who has ever thrown away a survey or refused to “answer a few questions” over the phone knows that getting people to participate in a survey isn’t easy. If the researcher wants to minimize bias, the best way to handle it is to get as many folks to respond as possible by following up, one, two, or even three times. Offer dollar bills, coupons, self-addressed stamped return envelopes, chances to win prizes, and so on. Every little bit helps.

If only those folks who feel very strongly respond to a survey, that means that only their opinions will count, because the other people who didn’t really care about the issue didn’t respond, and their “I don’t care” vote didn’t get counted. Or maybe they did care, but they just didn’t take the time to tell anyone. Either way, their vote doesn’t count.

For example, suppose 1,000 people are given a survey about whether the park rules should be changed to allow dogs without leashes. Most likely, the respondents would be those who strongly agree or disagree with the proposed rules. Suppose only 200 people responded — 100 against and 100 for the issue. That would mean that 800 opinions weren’t counted. Suppose none of those 800 people really cared about the issue either way. If you could count their opinions, the results would be $800 \div 1,000 = 80\%$ “no opinion,” $100 \div 1,000 = 10\%$ in favor of the new rules, and $100 \div 1,000 = 10\%$ against the new rules. But without the votes of the 800 non-respondents, the researchers would report, “Of the people who responded, 50% were in favor of the new rules and 50% were against them.” This gives the impression of a very different (and a very biased) result from the one you would’ve gotten if all 1,000 people had responded.

The *response rate* of a survey is a ratio found by taking the number of respondents

divided by the number of people who were originally asked to participate. You of course want to have the highest response rate you can get with your survey; but how high is high enough to be minimizing bias? The purest of the pure statisticians feel that a good response rate is anything over 70%, but I think we need to be a little more realistic. Today's fast-paced society is saturated with surveys; many if not most response rates fall far short of 70%. In fact, response rates for today's surveys are more likely to be in the 20% to 30% range, unless the survey is conducted by a professional polling organization such as Gallup or you are being offered a new car just for filling one out.



Look for the response rate when examining survey results. If the response rate is too low (much less than 50%) the results are likely to be biased and should be taken with a grain of salt, or even ignored.



Don't be fooled by a survey that claims to have a large number of respondents but actually has a low response rate; in this case, many people may have responded, but many more were asked and didn't respond.

Note that statistical formulas at this level (including the formulas in this book) assume that your sample size is equal to the number of respondents, so statisticians want you to know how important it is to follow up with people and not end up with biased data due to non-response. However, in reality, statisticians know that you can't always get everyone to respond, no matter how hard you try; indeed, even the U.S. Census doesn't have a 100% response rate. One way statisticians combat the non-response problem after the data have been collected is to break down the data to see how well it matches the target population. If it's a fairly good match, they can rest easier on the bias issue.

So which number do you put in for n in all those statistical formulas you use so often (such as the sample mean in Chapter 5)? You can't use the intended sample size (the number of people contacted). You have to use the final sample size (the number of people who responded). In the media you most often see only the number of respondents reported, but you also need the response rate (or the total number of respondents) to be able to critically evaluate the results.



Regarding the quality of results, selecting a smaller initial sample size and following them up more aggressively is a much better approach than selecting a larger group of potential respondents and having a low response rate, because of the bias introduced by non-response.

Interpreting results and finding problems

The purpose of a survey is to gain information about your target population; this

information can include opinions, demographic information, or lifestyles and behaviors. If the survey has been designed and conducted in a fair and accurate manner with the goals of the survey in mind, the data should provide good information as to what's happening with the target population (within the stated margin of error; see Chapter 12). The next steps are to organize the data to get a clear picture of what's happening; to analyze the data to look for links, differences, or other relationships of interest; and then to draw conclusions based on the results.

Organizing and analyzing

After a survey has been completed, the next step is to organize and analyze the data (in other words, crunch some numbers and make some graphs). Many different types of data displays and summary statistics can be created and calculated from survey data, depending on the type of information that was collected. (Numerical data, such as income, have different characteristics and are usually presented differently than categorical data, such as gender.) For more information on how data can be organized and summarized, see Chapters 5 through 7. Depending on the research question, different types of analyses can be performed on the data, including coming up with population estimates, testing a hypothesis about the population, or looking for relationships, to name a few. See Chapters 13, 14, 15, 18, and 19 for more on each of these analyses, respectively.



Watch for misleading graphs and statistics. Not all survey data are organized and analyzed fairly and correctly. See Chapter 3 for more about how statistics can go wrong.

Drawing conclusions

The conclusions are the best part of any survey — they're why the researchers do all of the work in the first place. If the survey was designed and carried out properly — the sample was selected carefully and the data were organized and summarized correctly — the results should fairly and accurately represent the reality of the target population. But, of course, not all surveys are done right. And even if a survey is done correctly, researchers can misinterpret or overinterpret results so that they say more than they really should.



You know the saying “Seeing is believing”? Some researchers are guilty of the converse, which is “Believing is seeing.” In other words, they claim to see what they want to believe about the results. All the more reason for you to know where the line is drawn between reasonable conclusions and misleading results, and to realize when others have crossed that line.

Here are some common errors made in drawing conclusions from surveys:

- ✓ Making projections to a larger population than the study actually represents
- ✓ Claiming a difference exists between two groups when a difference isn't really there (see Chapter 15)
- ✓ Saying, "these results aren't scientific, but . . .," and then going on to present the results as if they are scientific



To avoid common errors made when drawing conclusions, do the following:

1. Check whether the sample was selected properly and that the conclusions don't go beyond the population presented by that sample.

2. Look for any disclaimers about the survey *before* reading the results.

That way, if the results aren't based on a scientific survey (an accurate and unbiased survey), you'll be less likely to be influenced by the results you're reading. You can judge for yourself whether the survey results are credible.

3. Be on the lookout for statistically incorrect conclusions.

If someone reports a difference between two groups in terms of survey results, be sure that the difference is larger than the reported margin of error. If the difference is within the margin of error, you should expect the sample results to vary by that much just by chance, and the so-called "difference" can't really be generalized to the entire population. (See Chapter 14 for more on this.)



Know the limitations of any survey and be wary of any information coming from surveys in which those limitations aren't respected. A bad survey is cheap and easy to do, but you get what you pay for. But don't let big expensive surveys fool you either — they can be riddled with bias as well! Before looking at the results of any survey, investigate how it was designed and conducted, using the criteria and tips in this chapter, so you can judge the quality of the results and express yourself confidently and correctly about what is wrong.

Experiments: Medical Breakthroughs or Misleading Results?

In This Chapter

- ▶ Distinguishing experiments from observational studies
 - ▶ Dissecting the criteria for a good experiment
 - ▶ Watching for misleading results
-

Medical breakthroughs seem to come and go quickly. One day you hear about a promising new treatment for a disease, only to find out later that the drug didn't live up to expectations in the last stage of testing. Pharmaceutical companies bombard TV viewers with commercials for pills, sending millions of people to their doctors clamoring for the latest and greatest cures for their ills, sometimes without even knowing what the drugs are for. Anyone can search the Internet for details about any type of ailment, disease, or symptom and come up with tons of information and advice. But how much can you really believe? And how do you decide which options are best for you if you get sick, need surgery, or have an emergency?

In this chapter, you go behind the scenes of experiments, the driving force of medical studies and other investigations in which comparisons are made — comparisons that test, for example, which building materials are best, which soft drink teens prefer, and so on. You find out the difference between experiments and observational studies and discover what experiments can do for you, how they're supposed to be done, how they can go wrong, and how you can spot misleading results. With so many headlines, sound bites, and pieces of "expert advice" coming at you from all directions, you need to use all your critical thinking skills to evaluate the sometimes-conflicting information you're presented with on a regular basis.

Boiling Down the Basics of Studies

Although many different types of studies exist, you can basically boil them down to two types: experiments and observational studies. This section examines what exactly makes experiments different from other studies. But before I dive in to the details, I need to lay some jargon on you.

Looking at the lingo of studies

To understand studies, you need to find out what their commonly used terms mean:

- ✓ **Subjects:** Individuals participating in the study.
- ✓ **Observational study:** A study in which the researcher merely observes the subjects and records the information. No intervention takes place, no changes are introduced, and no restrictions or controls are imposed.
- ✓ **Experiment:** This study doesn't simply observe subjects in their natural state; it deliberately applies treatments to them in a controlled situation and studies their effects on the outcome.
- ✓ **Response:** The response is the variable whose outcome is the million dollar question; it's the variable whose outcome is of interest. For example, if researchers want to know what happens to your blood pressure when you take a large amount of Ibuprofen each day, the response variable is blood pressure.
- ✓ **Factor:** A factor is the variable whose effect on the response is being studied. For example, if you want to know whether a particular drug increases blood pressure, your factor is the amount of the drug taken. If you want to know which weight loss program is most effective, your factor would be the type of weight loss program used.

You can have more than one factor in a study; however, in this book I stick with discussing one factor only. For the analysis of two-factor studies, including the use of Analysis of Variance (ANOVA) and multiple comparisons to compare treatment combinations, you can check out my book *Statistics II For Dummies*, also published by Wiley.

- ✓ **Level:** A level is one possible outcome of a factor. Each factor has a certain number of levels. In the weight loss example, the factor is the type of weight loss program and the levels would be the specific programs studied (for example Weight Watchers, South Beach, or the famous Potato Diet). Levels need not be ascending in any way; however, in a study like the drug example, the levels would be the various dosages taken each day, in increasing amounts.
- ✓ **Treatment:** A treatment is a combination of the levels of the factors being studied. If you only have one factor, the levels and the treatments are the same thing. If you have more than one factor, each combination of levels of the factors is called a treatment.

For example, if you want to study the effects of the type of weight loss program and the amount of water consumed daily, you have two factors: 1) the type of program, with 3 levels (Weight Watchers, South Beach, Potato Diet); and 2) the amount of water consumed, with, say, 3 levels (24, 48, and 64 ounces per day). In

this case, there are $3 * 3 = 9$ treatments: Weight Watchers and 24 ounces of water per day; Weight Watchers and 48 ounces of water per day, . . . all the way up to the famous Potato Diet and 64 ounces of water per day. Each subject is assigned to one treatment. (With my luck, I'd get that last treatment.)

- ✓ **Cause and effect:** A factor and a response have a cause-and-effect relationship if a change in the factor results in a direct change in the response (for example, increasing calorie intake causes weight gain).

In the following sections, you see the differences between observational studies and experiments, when each is used, and what their strengths and/or weaknesses may be.

Observing observational studies

Just like with tools, you want to find the right type of study for the right job. In certain situations, observational studies are the optimal way to go. The most common observational studies are *polls* and *surveys* (see Chapter 16). When the goal is simply to find out what people think and to collect some demographic information (such as gender, age, income, and so on), surveys and polls can't be beat, as long as they're designed and conducted correctly.

In other situations, especially those looking for cause-and-effect relationships, observational studies aren't optimal. For example, suppose you took a couple of vitamin C pills last week; is that what helped you avoid getting that cold that's going around the office? Maybe the extra sleep you got recently or the extra hand-washing you've been doing helped you ward off the cold. Or maybe you just got lucky this time. With so many variables in the mix, how can you tell which one had an influence on the outcome of your not getting a cold? An experiment that takes these other variables into account is the way to go.



When looking at the results of any study, first determine what the purpose of the study was and whether the type of study fits the purpose. For example, if an observational study was done instead of an experiment to establish a cause-and-effect relationship, any conclusions that are drawn should be carefully scrutinized.

Examining experiments

The object of an experiment is to see if the response changes as a result of the factor you are studying; that is, you are looking for cause and effect. For example, does taking Ibuprofen cause blood pressure to increase? If so, by how much? But because results will vary with any experiment, you want to know that your results have a high chance of being repeatable if you found something interesting happening. That is, you want to know that your results were unlikely to be due to chance; statisticians call such results



A good experiment is conducted by creating a very controlled environment — so controlled that the researcher can pinpoint whether a certain factor or combination of factors causes a change in the response variable, and if so, the extent to which that factor (or combination of factors) influences the response. For example, to gain government approval for a proposed blood pressure drug, pharmaceutical researchers set up experiments to determine whether that drug helps lower blood pressure, what dosage level is most appropriate for each different population of patients, what side effects (if any) occur, and to what extent those side effects occur in each population.

Designing a Good Experiment

How an experiment is designed can mean the difference between good results and garbage. Because most researchers are going to write the most glowing press releases that they can about their experiments, you have to be able to sort through the hype to determine whether to believe the results you're being told. To decide whether an experiment is credible, check to see if it meets *all* the following criteria for a good experiment. A good experiment:

- ✓ **Makes comparisons**
- ✓ **Includes a large enough sample size so that the results are accurate**
- ✓ **Chooses subjects that most accurately represent the target population**
- ✓ **Assigns subjects randomly to the treatment group(s) and the control group**
- ✓ **Controls for possible confounding variables**
- ✓ **Is ethical**
- ✓ **Collects good data**
- ✓ **Applies the proper data analysis**
- ✓ **Makes appropriate conclusions**

In this section, each of these criteria is explained and illustrated with examples.

Designing the experiment to make comparisons

Every experiment has to make bonafide comparisons to be credible. This seems to go

without saying, but researchers often are so gung-ho to prove their results that they forget (or just don't bother) to show that their factor, and not some other factor(s), including random chance, was the actual cause for any differences found in the response.

For example, suppose a researcher is convinced that taking vitamin C prevents colds, and she assigns subjects to take one vitamin C pill per day and follows them for 6 months. Suppose the subjects get very few colds during that time. Can she attribute these results to the vitamin C and nothing else? No; there's no way of knowing whether the subjects would have been just as healthy without the vitamin C, due to some other factor(s), or just by chance. There's nothing to compare the results to.



To tease out the real effect (if any) that your factor has on the response, you need a baseline to compare the results to. This baseline is called the *control*. Different methods exist for creating a control in an experiment; depending on the situation, one method typically rises to the top as being the most appropriate. Three common methods for including control are to administer: 1) a fake treatment; 2) a standard treatment; or 3) no treatment. The following sections describe each method.



When examining the results of an experiment, make sure the researchers established a baseline by creating a control group. Without a control group, you have nothing to compare the results to, and you never know whether the treatment being applied was the real cause of any differences found in the response.

Fake treatments — the placebo effect

A fake treatment (also called a *placebo*) is not distinguishable from a “real” treatment by the subject. For example, when drugs are administered, a subject assigned to the placebo will receive a fake pill that looks and tastes exactly like a real pill; it's just filled with an inert substance like sugar instead of the actual drug. A placebo establishes a baseline measure for what responses would have taken place anyway, in lieu of any treatment (this would have helped the vitamin C study mentioned under “Designing the experiment to make comparisons”). But a fake treatment also takes into account what researchers call the *placebo effect*, a response that people have (or think they're having) because they know they're getting some type of “treatment” (even if that treatment is a fake treatment, such as sugar pills).

Pharmaceutical companies are required to account for the placebo effect when examining both the positive and negative effects of a drug. When you see an ad for a drug in a magazine, you see the positive results of the drug standing out in big, bright, happy, colorful visuals. Then look at the back of the page and you see it's entirely filled in black with words written in 3-point font. Embedded somewhere on that page, you can find one or more tiny tables that show the number and nature of side effects reported by each *treatment group* (subjects who received an actual treatment) as well as the *control group*.

(subjects who were administered a placebo).



If the control group is on a placebo, you may expect the subjects not to report any side effects, but you would be wrong. If you are taking a pill, you know it could be an actual drug, and you are being asked whether or not you're experiencing side effects, you might be surprised at what your response would be.

If you don't take the placebo effect into account, you have to believe that any side effects (or positive results) reported are actually due to the drug. This gives an artificially high number of reported side effects because at least some of those reports are likely due to the placebo effect and not to the drug itself. If you have a control group to compare with, you can subtract the percentage of people in the control group who reported the side effects from the percentage of people in the treatment group that reported the side effects, and examine the magnitude of the numbers that remain. You're in essence looking at the net number of reported side effects due to the drug, rather than the gross number of side effects, some of which are due to the placebo effect.



The placebo effect has been shown to be real. If you want to be fair about examining the reported side effects (or positive reactions) of a treatment, you have to also take into account the side effects (or positive reactions) that the control group reports — those reactions that are due to the placebo effect only.

Standard treatments



In some situations, such as when the subjects have very serious diseases, offering a fake treatment as an option may be unethical. One famous example of a breech in ethics occurred in 1997. The U.S. government was harshly criticized for financing an HIV study that examined new dosage levels of AZT, a drug known at that time to cut the risk of HIV transmission from pregnant mothers to their babies by two-thirds. This particular study, in which 12,000 pregnant women with HIV in Africa, Thailand, and the Dominican Republic participated, had a deadly design. Researchers gave half of the women various dosages of AZT, but the other half of the women received sugar pills. Of course, had the U.S. government realized that a placebo was being given to half of the subjects, it wouldn't have supported the HIV study. It's not ethical to give a fake treatment to anyone with a deadly disease for which a standard treatment is available (in this case, the standard dosage of AZT).

When ethical reasons bar the use of fake treatments, the new treatment is compared to at least one existing or standard treatment that is known to be an effective treatment. After researchers have enough data to see that one of the treatments is working better than the other, they generally stop the experiment and put everyone on the better treatment;

again, for ethical reasons.

No treatment

“No treatment” means the researcher can’t help but tell which group the subject is in, due to the nature of the experiment. The subjects in this case aren’t receiving any type of intervention in terms of their behavior, but they still serve as a control, establishing a baseline of data to compare their results with those in the treatment group(s). For example, if you want to determine whether speed walking around the block ten times a day lowers a person’s resting heart rate after six months, the subjects in your control group know they aren’t going to be speed walking — obviously you can’t do fake speed walking (although faking exercising and still reaping the benefits would be great, wouldn’t it?).



In situations where the control group receives no treatment, you still make sure the groups of subjects (speed walkers versus non-speed walkers) are similar in as many ways as possible, and that the other criteria for a good experiment are being met. (See “Designing a Good Experiment” for the list of criteria.)

Selecting the sample size

The size of a (good) sample greatly affects the accuracy of the results. The larger the sample size, the more accurate the results, and the more powerful the statistical tests (in terms of being able to detect real results when they exist). In this section, I hit the highlights; Chapter 14 has the details.



The word *sample* is often attributed to surveys where a random sample is selected from the target population (see Chapter 16). However, in the setting of experiments, a sample means the group of subjects who have volunteered to participate.

Limiting small samples to small conclusions

You may be surprised at the number of research headlines that have been made regarding large populations that were based on very small samples. Such headlines can be of concern to statisticians, who know that detecting true statistically significant results in a large population using a small sample is difficult because small data sets have more variability from sample to sample (see Chapter 12). When sample sizes are small and big conclusions have been made by the researcher, either the researchers didn’t use the right hypothesis test to analyze their data (for example, using the Z-distribution rather than the *t*-distribution; see Chapter 10) or the difference was so large that it would be very difficult to miss. The latter isn’t always the case, however.



Be wary of research conclusions that find significant results based on small sample sizes (especially for experiments involving many treatments but only a few subjects assigned to each treatment). Statisticians want to see at least five subjects per treatment, but (much) more is (much) better. You do need to be aware of some of the limitations of experiments such as cost, time, as well as ethical issues, and realize that the number of subjects for experiments is often smaller than the number of participants in a survey.

If the results are important to you, ask for a copy of the research report and look to see what type of analysis was done on the data. Also look at the sample of subjects to see whether this sample truly represents the population about which the researchers are drawing conclusions.

Defining sample size

When asking questions about *sample size*, be specific about what you mean by the term. For example, you can ask how many subjects were selected to participate and also ask for the number who actually completed the experiment; these two numbers can be very different. Make sure the researchers can explain any situations in which the research subjects decided to drop out or were unable (for some reason) to finish the experiment.

For example, an article in *The New York Times* titled “Marijuana Is Called an Effective Relief in Cancer Therapy” says in the opening paragraph that marijuana is “far more effective” than any other drug in relieving the side effects of chemotherapy. When you get into the details, you find out that the results are based on only 29 patients (15 on the treatment, 14 on a placebo). Then you find out that only 12 of the 15 patients in the treatment group actually completed the study. What happened to the other three subjects?



Sometimes researchers draw their conclusions based on only those subjects who completed the study. This can be misleading, because the data don’t include information about those who dropped out (and why), which may be leading to biased data. For a discussion of the sample size you need to achieve a certain level of accuracy, see Chapter 13.



Accuracy isn’t the only issue in terms of having “good” data. You still need to worry about eliminating bias by selecting a random sample (see Chapter 16 for more on how random samples are taken).

Choosing the subjects

The first step in carrying out an experiment is selecting the subjects (participants). Although researchers would like their subjects to be selected randomly from their respective populations, in most cases, this just isn't appropriate. For example, suppose a group of eye researchers wants to test out a new laser surgery on nearsighted people. They need a random sample of subjects, so they randomly select various eye doctors from across the country and randomly select nearsighted patients from these doctors' files. They call up each person selected and say, "We're experimenting with a new laser surgery technique for nearsightedness, and you've been selected at random to participate in our study. When can you come in for the surgery?" Something tells me that this approach wouldn't go over very well with many people receiving the call (although some would probably jump at the chance, especially if they didn't have to pay for the procedure).



The point is that getting a truly random sample of people to participate in an experiment is generally more difficult than getting a random sample of folks to participate in a survey. However, statisticians can build techniques into the design of an experiment to help minimize the potential bias that can occur.

Making random assignments

One way to minimize bias in an experiment is to introduce some randomness. After the sample has been decided on, the subjects are randomly divided into treatment and control groups. The treatment groups receive the various treatments being studied, and the control group receives the current (or standard) treatment, no treatment, or a placebo. (See the section "Designing the experiment to make comparisons" earlier in this chapter.)

Making random assignments of subjects to treatments is an extremely critical step toward minimizing bias in an experiment. Suppose a researcher wants to determine the effects of exercise on heart rate. The subjects in his treatment group run 5 miles and have their heart rates measured before and after the run. The subjects in his control group sit on the couch the whole time and watch reruns of old TV shows. Which group would you rather be in? Some health nuts out there would no doubt volunteer for the treatment group. If you're not crazy about the idea of running five miles, you may opt for the easy way out and volunteer to be a couch potato. (Or maybe you hate to watch old reruns so much that you'd run five miles to avoid that.)

Finding volunteers

To find subjects for their experiments, researchers often advertise for volunteers and offer them incentives such as money, free treatments, or follow-up care for their participation. Medical research on humans is complicated and difficult, but it's necessary in order to really know whether a treatment works, how well it works, what the dosage should be, and what the side effects are. In order to

prescribe the right treatments in the right amounts in real-life situations, doctors and patients depend on these studies being representative of the general population. In order to recruit such representative subjects, researchers have to do a broad advertisement campaign and select enough participants with enough different characteristics to represent a cross section of the populations of folks who will be prescribed these treatments in the future.

What impact would this selective volunteering have on the results of the study? If only the health nuts (who probably already have excellent heart rates) volunteer to be in the treatment group, the researcher will be looking only at the effect of the treatment (running five miles) on very healthy and active people. He won't see the effect that running five miles has on the heart rates of couch potatoes. This non-random assignment of subjects to the treatment and control groups could have a huge impact on the conclusions he draws from this study.



To avoid major bias in the results of an experiment, subjects must be randomly assigned to treatments by a third party and not be allowed to choose which group they will be in. The goal of random assignment is to create homogenous groups; any unusual characteristics or biases have an equal chance of appearing in any of the groups. Keep this in mind when you evaluate the results of an experiment.

Controlling for confounding variables

Suppose you're participating in a research study that looks at factors influencing whether you catch a cold. If a researcher records only whether you got a cold after a certain period of time and asks questions about your behavior (how many times per day you washed your hands, how many hours of sleep you get each night, and so on), the researcher is conducting an observational study. The problem with this type of observational study is that without controlling for other factors that may have had an influence and without regulating which action you were taking when, the researcher won't be able to single out exactly which of your actions (if any) actually impacted the outcome.



The biggest limitation of observational studies is that they can't really show true cause-and-effect relationships, due to what statisticians call confounding variables. A *confounding variable* is a variable or factor that was not controlled for in the study but can have an influence on the results.

For example, one news headline boasted, "Study links older mothers, long life." The opening paragraph said that women who have a first baby after age 40 have a much better chance of living to be 100, compared to women who have a first baby at an earlier age. When you get into the details of the study (done in 1996) you find out, first of all, that it

was based on 78 women in suburban Boston who were born in 1896 and had lived to be at least 100, compared to 54 women who were also born in 1896 but died in 1969 (the earliest year the researchers could get computerized death records). This so-called “control group” lived to be exactly 73, no more and no less. Of the women who lived to be at least 100 years of age, 19% had given birth after age 40, whereas only 5.5% of the women who died at age 73 had given birth after age 40.

I have a real problem with these conclusions. What about the fact that the “control group” was based only on mothers who died in 1969 at age 73? What about all the other mothers who died *before* age 73, or who died between the ages of 73 and 100? What about other variables that may affect both mothers’ ages at the births of their children and longer life spans — variables such as financial status, marital stability, or other socioeconomic factors? The women in this study were in their thirties during the Depression; this may have influenced both their life span and if or when they had children.



How do researchers handle confounding variables? They control for them as best they can, for as many of them as they can anticipate, trying to minimize their possible effect on the response. In experiments involving human subjects, researchers have to battle many confounding variables.

For example, in a study trying to determine the effect of different types and volumes of music on the amount of time grocery shoppers spend in the store (yes, they do think about that), researchers have to anticipate as many possible confounding variables ahead of time and then control for them. What other factors besides volume and type of music may influence the amount of time you spend in a grocery store? I can think of several factors: gender, age, time of day, whether you have children with you, how much money you have, the day of the week, how clean and inviting the store is, how nice the employees are, and (most importantly) what your motive is — are you shopping for the whole week, or are you just running in to grab a candy bar?

How can researchers begin to control for so many possible confounding factors? Some of them can be controlled for in the design of the study, such as the time of the day, day of the week, and reason for shopping. But other factors (such as the perception of the store environment) depend totally on the individual in the study. The ultimate form of control for those person-specific confounding variables is to use pairs of people that are matched according to important variables, or to just use the same person twice: once with the treatment and once without. This type of experiment is called a *matched-pairs design*. (See Chapter 15 for more on this.)



Before believing any medical headlines (or any headlines with statistics, for that matter), look to see how the study was conducted. Observational studies can’t control for confounding variables, so their results are not as statistically meaningful

(no matter what the statistics say) as the results of a well-designed experiment are. In cases where an experiment can't be done (after all, no one can force you to have a baby after or before age 40), make sure the observational study is based on a large enough sample that represents a cross-section of the population. And think about possible confounding variables that may affect the conclusions being drawn.

Respecting ethical issues

The trouble with experiments is that some experimental designs are not ethical. You can't force research subjects to smoke in order to see whether they get lung cancer, for example — you can only look at people who have lung cancer and work backward to see what *factors* (variables being studied) may have caused the disease. But because you can't control for the various factors you're interested in — or for any other variables, for that matter — singling out any one particular cause becomes difficult with observational studies. That's why so much evidence was needed to show that smoking causes lung cancer, and why the tobacco companies only recently had to pay huge penalties to victims.

Although the causes of cancer and other diseases can't be determined ethically by conducting experiments on humans, new treatments for cancer can be (and are) tested using experiments. Medical studies that involve experiments are called *clinical trials*. The U.S. government has a registry of federally and privately supported clinical trials conducted in the United States and around the world; it also has information available on who may participate in various clinical trials. Check out www.clinicaltrials.gov for more information.

Serious experiments (such as those funded by and/or regulated by the U.S. government) must pass a huge series of tests that can take years to carry out. The approval of a new drug, for example, goes through a very lengthy, comprehensive, and detailed process regulated and monitored by the FDA (Federal Drug Administration). One reason the cost of prescription drugs is so high is the massive amount of time and money needed to conduct research and development of new drugs, most of which fail to pass the tests and have to be scrapped.

Any experiments involving human subjects are also regulated by the federal government and have to gain approval by a committee created for the purpose of protecting "the rights and welfare of the participants." The committees set up for different organizations have different names (such as Institutional Review Board [IRB], Independent Ethics Committee [IEC], or Ethical Review Board [ERB], to name a few) but they all serve the same purpose. Research conducted on animals is more nebulous in terms of regulations and continues to be a topic of much debate and controversy in the U.S. and around the world.



Surveys, polls, and other observational studies are fine if you want to know people's opinions, examine their lifestyles without intervention, or examine some demographic variables. If you want to try to determine the cause of a certain outcome or behavior (that is, a reason why something happened), an experiment is a much better way to go. If an experiment isn't possible because of ethics concerns (or because of expense or other reasons), a large body of observational studies examining many different factors and coming up with similar conclusions is the next best thing. (See Chapter 18 for more about cause-and-effect relationships.)

Collecting good data

What constitutes “good” data? Statisticians use three criteria for evaluating data quality; each of the criteria really relates most strongly to the quality of the measurement instrument that’s used in the process of collecting the data. To decide whether you’re looking at good data from a study, look for these characteristics:

- ✓ **The data are reliable — you can get repeatable results with subsequent measurements.** Many bathroom scales give unreliable data. You get on the scale, and it gives you one number. You don’t believe the number, so you get off, get back on, and get a different number. (If the second number is lower, you’ll most likely quit at this point; if not, you may continue getting on and off until you see a number you like.) Or you can do what some researchers do: Take three measurements, find the average, and use that; at least this will improve the reliability a bit.

Unreliable data come from unreliable measurement instruments or unreliable data collection methods. Errors can go beyond the actual scales to more intangible measurement instruments, like survey questions, which can give unreliable results if they’re written in an ambiguous way (see Chapter 16).



Find out how the data were collected when examining the results of a study. If the measurements are unreliable, the data could be inaccurate.

- ✓ **The data are valid — they measure what they’re supposed to measure.** Checking the validity of data requires you to step back and look at the big picture. You have to ask the question: Do these data measure what they should be measuring? Or should the researchers have been collecting altogether different data? The appropriateness of the measurement instrument used is important. For example, many educators say that a student’s transcript is not a valid measure of their ability to perform well in college. Alternatives include a more holistic approach, taking into account not only grades, but adding weight to elements such as service, creativity, social involvement, extracurricular activities, and the like.



Before accepting the results of an experiment, find out what data were measured and how they were measured. Be sure the researchers are collecting valid data that are appropriate for the goals of the study.

- ✓ **The data are unbiased — they contain no systematic errors that either add to or subtract from the true values.** Biased data are data that systematically overmeasure or undermeasure the true result. Bias can occur almost anywhere during the design or implementation of a study. Bias can be caused by a bad measurement instrument (like that bathroom scale that's "always" 5 pounds over), by survey questions that lead participants in a certain way, or by researchers who know what treatment each subject received and who have preconceived expectations.



Bias is probably the number-one problem in collecting good data. However, you can minimize bias with methods similar to those discussed in Chapter 16 for surveys and in the "Making random assignments" section earlier in this chapter, and by making your experiments double-blind whenever possible.

Double-blind means neither the subjects nor the researchers know who got what treatment or who is in the control group. The subjects need to be oblivious to which treatment they're getting so that the researchers can measure the placebo effect. And researchers should be kept in the dark so they don't treat subjects differently by either expecting or not expecting certain responses from certain groups. For example, if a researcher knows you're in the treatment group to study the side effects of a new drug, she may expect you to get sick and therefore may pay more attention to you than if she knew you were in the control group. This can result in biased data and misleading results.

If the researcher knows who got what treatment but the subjects don't know, the study is called a *blind* study (rather than a double-blind study). Blind studies are better than nothing, but double-blind studies are best. In case you're wondering: In a double-blind study, does *anyone* know which treatment was given to which subjects? Relax; typically a third party, such as a lab assistant, does that part.

In some cases the subjects know which group they're in because it's unconcealable — for example, when comparing the benefits of doing yoga versus jogging. However, bias can be reduced by not telling the subjects the precise purpose of the study. This irregular type of plan would have to be reviewed by an institutional review board to make sure it isn't unethical to do; see the earlier section "Respecting ethical issues."

Analyzing the data properly

After the data have been collected, they're put into that mysterious box called the *statistical analysis for number crunching*. The choice of analysis is just as important (in terms of the quality of the results) as any other aspect of a study. A proper analysis should be planned in advance, during the design phase of the experiment. That way, after the data are collected, you won't run into any major problems during the analysis.

Here's the bottom line when selecting the proper analysis: Ask yourself the question, "After the data are analyzed, will I be able to legitimately and correctly answer the question that I set out to answer?" If the answer is "no," then that analysis isn't appropriate.

Some basic types of statistical analyses include *confidence intervals* (used when you're trying to estimate a population value, or the difference between two population values); *hypothesis tests* (used when you want to test a claim about one or two populations, such as the claim that one drug is more effective than another); and *correlation and regression analyses* (used when you want to show if and/or how one quantitative variable can predict or cause changes in another quantitative variable). See Chapters 13, 15, and 18, respectively, for more on each of these types of analyses.



When choosing how you're going to analyze your data, you have to make sure that the data and your analysis will be compatible. For example, if you want to compare a treatment group to a control group in terms of the amount of weight lost on a new (versus an existing) diet program, you need to collect data on how much weight each person lost — not just each person's weight at the end of the study.

Making appropriate conclusions

In my opinion, the biggest mistakes researchers make when drawing conclusions about their studies are the following (discussed in the following sections):

- ✓ Overstating their results
- ✓ Making connections or giving explanations that aren't backed up by the statistics
- ✓ Going beyond the scope of the study in terms of whom the results apply to

Overstating the results

Many times, the headlines in the media overstate actual research results. When you read a headline or otherwise hear about a study, be sure to look further to find out the details of how the study was done and exactly what the conclusions were.

Press releases often overstate results, too. For example, in a recent press release by the National Institute for Drug Abuse, the researchers claimed that use of the street drug

Ecstasy was down from the previous year. However, when you look at the actual statistical results in the report, you find that the percentage of teens *from the sample* who said they'd used Ecstasy was lower than those from the previous year, but this difference was not found to be statistically significant when they tried to project it onto the population of *all* teens. This discrepancy means that although fewer teens in the sample used Ecstasy that year, the difference wasn't enough to account for more than chance variability from sample to sample. (See Chapter 14 for more about statistical significance.)



Headlines and leading paragraphs in press releases and news articles often overstate the actual results of a study. Big results, spectacular findings, and major breakthroughs make the news these days, and reporters and others in the media constantly push the envelope in terms of what is and isn't newsworthy. How can you sort out the truth from exaggeration? The best thing to do is to read the fine print.

Taking the results one step beyond the actual data

A study that links having children later in life to longer life spans illustrates another point about research results. Do the results of this observational study mean that having a baby later in life can make you live longer? "No," said the researchers. Their explanation of the results was that having a baby later in life may be due to women having a "slower" biological clock, which presumably would then result in the aging process being slowed down.

My question to these researchers is, "Then why didn't you study *that*, instead of just looking at their ages?" The study didn't include any information that would lead me to conclude that women who had children after age 40 aged at a slower rate than other women, so in my view, the researchers shouldn't make that conclusion. Or the researchers should state clearly that this view is only a theory and requires further study. Based on the data in this study, the researchers' theory seems like a leap of faith (although since I became a new mom at age 41, I'll hope for the best!).

Frequently in a press release or news article, the researcher will give an explanation about *why* he thinks the results of the study turned out the way they did and what implications these results have for society as a whole when the "why" hasn't been studied yet. These explanations may have been in response to a reporter's questions about the research — questions that were later edited out of the story, leaving only the juicy quotes from the researcher. Many of these after-the-fact explanations are no more than theories that have yet to be tested. In such cases, you should be wary of conclusions, explanations, or links drawn by researchers that aren't backed up by their studies.



Be aware that the media wants to make you read the article (they get paid to do that), so they will have strong headlines, or will make unconfirmed “cause-effect” statements because it is their job to sell the story. It is *your* job to be wary.

Generalizing results to people beyond the scope of the study

You can make conclusions only about the population that's represented by your sample. If you sample men only, you can't make conclusions about women. If you sample healthy young people, you can't make your conclusions about everyone. But many researchers try to do just that, and it can give misleading results.

Here's how you can determine whether a researcher's conclusions measure up (Chapter 16 has more on samples and populations):

- 1. Find out what the target population is (that is, the group that the researcher wants to make conclusions about).**
- 2. Find out how the sample was selected and see whether the sample is representative of that target population (and not some more narrowly defined population).**
- 3. Check the conclusions made by the researchers and make sure they're not trying to apply their results to a broader population than they actually studied.**

Making Informed Decisions

Just because someone says they conducted a “scientific study” or a “scientific experiment” doesn’t mean it was done right or that the results are credible (not that I’m saying you should discount everything that you see and hear). Unfortunately, I’ve come across a lot of bad experiments in my days as a statistical consultant. The worst part is that if an experiment was done poorly, you can’t do anything about it after the fact except ignore the results — and that’s exactly what you need to do.



Here are some tips that help you make an informed decision about whether to believe the results of an experiment, especially one whose results are very important to you:

- ✓ **When you first hear or see the result, grab a pencil and write down as much as you can about what you heard or read, where you heard or read it, who did the research, and what the main results were.** (I keep pencil and paper in my TV room and in my purse just for this purpose.)
- ✓ **Follow up on your sources until you find the person who did the original**

research and then ask them for a copy of the report or paper.

- ✓ **Go through the report and evaluate the experiment according to the eight steps for a good experiment described in the “Designing a Good Experiment” section of this chapter.** (You really don’t have to understand everything written in a report in order to do that.)
- ✓ **Carefully scrutinize the conclusions that the researcher makes regarding his or her findings.** Many researchers tend to overstate results, make conclusions beyond the statistical evidence, or try to apply their results to a broader population than the one they studied.
- ✓ **Never be afraid to ask questions of the media, the researchers, and even your own experts.** For example, if you have a question about a medical study, ask your doctor. He or she will be glad that you’re an empowered and well-informed patient!
- ✓ **And finally, don’t get overly skeptical, just because you’re now a lot more aware of all the bad practices going on out there.** Not everything is bad. There are many more good researchers, credible results, and well-informed reporters than not. You have to maintain a sense of being cautious and ready to spot problems without discounting everything.

Chapter 18

Looking for Links: Correlation and Regression

In This Chapter

- ▶ Exploring statistical relationships between numerical variables
 - ▶ Looking at correlation and linear regression
 - ▶ Making predictions based on known relationships
 - ▶ Considering correlation versus causation
-

Today's media provide a steady stream of information, including reports on all the latest links that have been found by researchers. Just today I heard that increased video game use can negatively affect a child's attention span, the amount of a certain hormone in a woman's body can predict when she will enter menopause, and the more depressed you get, the more chocolate you eat, and the more chocolate you eat, the more depressed you get (how depressing!).

Some studies are truly legitimate and help improve the quality and longevity of our lives. Other studies are not so clear. For example, one study says that exercising 20 minutes three times a week is better than exercising 60 minutes one time a week, another study says the opposite, and yet another study says there is no difference.

If you are a confused consumer when it comes to links and correlations, take heart; this chapter can help. You'll gain the skills to dissect and evaluate research claims and make your own decisions about those headlines and sound bites that you hear each day alerting you to the latest correlation. You'll discover what it truly means for two variables to be correlated, when a cause-and-effect relationship can be concluded, and when and how to predict one variable based on another.

Picturing a Relationship with a Scatterplot

An article in *Garden Gate* magazine caught my eye: "Count Cricket Chirps to Gauge Temperature." According to the article, all you have to do is find a cricket, count the number of times it chirps in 15 seconds, add 40, and voilà! You've just estimated the temperature in Fahrenheit.

The National Weather Service Forecast Office even puts out its own "Cricket Chirp Converter." You enter the number of cricket chirps recorded in 15 seconds, and the converter gives you the estimated temperature in four different units, including

Fahrenheit and Celsius.

A fair amount of research does support the claim that frequency of cricket chirps is related to temperature. For the purpose of illustration I've taken only a subset of some of the data (see Table 18-1).

Table 18-1 Cricket Chirps and Temperature Data (Excerpt)

<i>Number of Chirps (in 15 Seconds)</i>	<i>Temperature (Fahrenheit)</i>
18	57
20	60
21	64
23	65
27	68
30	71
34	74
39	77

Notice that each observation is composed of two variables that are tied together: the number of times the cricket chirped in 15 seconds (the *X*-variable) and the temperature at the time the data was collected (the *Y*-variable). Statisticians call this type of two-dimensional data *bivariate* data. Each observation contains one pair of data collected simultaneously. For example, row one of Table 18-1 depicts a pair of data (18, 57).

Bivariate data is typically organized in a graph that statisticians call a *scatterplot*. A scatterplot has two dimensions, a horizontal dimension (the *X*-axis) and a vertical dimension (the *Y*-axis). Both axes are numerical; each one contains a number line. In the following sections, I explain how to make and interpret a scatterplot.

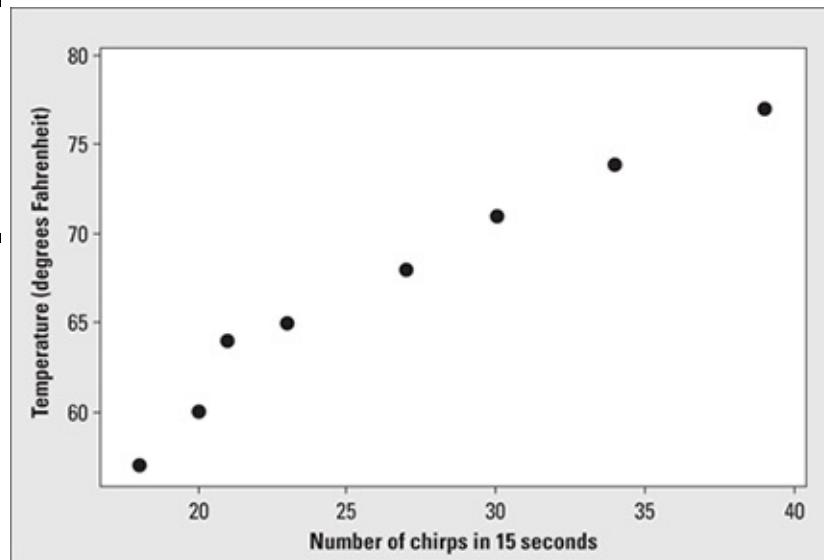
Making a scatterplot



Placing observations (or points) on a scatterplot is similar to playing the game Battleship. Each observation has two coordinates; the first corresponds to the first piece of data in the pair (that's the *X* coordinate; the amount that you go left or right). The second coordinate corresponds to the second piece of data in the pair (that's the *Y*-coordinate; the amount that you go up or down). You place the point representing that observation at the intersection of the two coordinates.

Figure 18-1 shows a scatterplot for the cricket chirps and temperature data listed in Table 18-1. Because I ordered the data according to their *X*-values, the points on the scatterplot correspond from left to right to the observations given in Table 18-1, in the order listed.

Figure 18-1:
Scatterplot of cricket chirps in relation to outdoor temperature.



Interpreting a scatterplot



You interpret a scatterplot by looking for trends in the data as you go from left to right:

- ✓ If the data show an uphill pattern as you move from left to right, this indicates a *positive relationship between X and Y*. As the X-values increase (move right), the Y-values increase (move up) a certain amount.
- ✓ If the data show a downhill pattern as you move from left to right, this indicates a *negative relationship between X and Y*. As the X-values increase (move right) the Y-values decrease (move down) by a certain amount.
- ✓ If the data don't seem to resemble any kind of pattern (even a vague one), then no relationship exists between X and Y.

One pattern of special interest is a *linear* pattern, where the data has a general look of a line going uphill or downhill. Looking at Figure 18-1, you can see that a positive linear relationship does appear between number of cricket chirps and the temperature. That is, as the cricket chirps increase, the temperature increases as well.



In this chapter I explore linear relationships only. A *linear relationship between X and Y* exists when the pattern of X- and Y-values resembles a line, either uphill (with a positive slope) or downhill (with a negative slope). Other types of trends may exist in addition to the uphill/downhill linear trends (for example, curves or exponential functions); however, these trends are beyond the scope of this book. The good news is that many relationships do fall under the uphill/downhill linear scenario.



Scatterplots show possible associations or relationships between two variables. However, just because your graph or chart shows something is going on, it doesn't mean that a cause-and-effect relationship exists.

For example, a doctor observes that people who take vitamin C each day seem to have fewer colds. Does this mean vitamin C prevents colds? Not necessarily. It could be that people who are more health conscious take vitamin C each day, but they also eat healthier, are not overweight, exercise every day, and wash their hands more often. If this doctor really wants to know if it's the vitamin C that's doing it, she needs a well-designed experiment that rules out these other factors. (See the later section "Explaining the Relationship: Correlation versus Cause and Effect" for more information.)

Quantifying Linear Relationships Using the Correlation

After the bivariate data have been organized graphically with a scatterplot (see the preceding section), and you see some type of linear pattern, the next step is to do some statistics that can quantify or measure the extent and nature of the relationship. In the following sections, I discuss *correlation*, a statistic measuring the strength and direction of a linear relationship between two variables; in particular, how to calculate and interpret correlation and understand its most important properties.

Calculating the correlation

In the earlier section "Interpreting a scatterplot," I say data that resembles an uphill line has a positive linear relationship and data that resembles a downhill line has a negative linear relationship. However, I didn't address the issue of whether or not the linear relationship was strong or weak. The strength of a linear relationship depends on how closely the data resembles a line, and of course varying levels of "closeness to a line" exist.

Can one statistic measure both the strength and direction of a linear relationship between two variables? Sure! Statisticians use the *correlation coefficient* to measure the strength and direction of the linear relationship between two numerical variables X and Y . The correlation coefficient for a sample of data is denoted by r .



Although the street definition of *correlation* applies to any two items that are related (such as gender and political affiliation), statisticians use this term only in the context of two numerical variables. The formal term for correlation is the

correlation coefficient. Many different correlation measures have been created; the one used in this case is called the *Pearson correlation coefficient* (but from now on I'll just call it the correlation).

The formula for the correlation (r) is

$$r = \frac{1}{n-1} \left(\frac{\sum_{x,y} (x-\bar{x})(y-\bar{y})}{s_x s_y} \right)$$

where n is the number of pairs of data; \bar{x} and \bar{y} are the sample means of all the x -values and all the y -values, respectively; and s_x and s_y are the sample standard deviations of all the x - and y -values, respectively.



Use the following steps to calculate the correlation, r , from a data set:

1. Find the mean of all the x -values (\bar{x}) and the mean of all the y -values (\bar{y}).

See Chapter 5 for more on calculating the mean.

2. Find the standard deviation of all the x -values (call it s_x) and the standard deviation of all the y -values (call it s_y).

See Chapter 5 to find out how to calculate the standard deviation.

3. For each (x, y) pair in the data set, take x minus \bar{x} and y minus \bar{y} , and multiply them together to get $(x-\bar{x})(y-\bar{y})$.

4. Add up all the results from Step 3.

5. Divide the sum by $s_x * s_y$.

6. Divide the result by $n - 1$, where n is the number of (x, y) pairs. (It's the same as multiplying by 1 over $n - 1$.)

This gives you the correlation, r .

For example, suppose you have the data set $(3, 2)$, $(3, 3)$, and $(6, 4)$. You calculate the correlation coefficient r via the following steps. (Note for this data the x -values are 3, 3, 6, and the y -values are 2, 3, 4.)

1. \bar{x} is $12 \div 3 = 4$, and \bar{y} is $9 \div 3 = 3$.

2. The standard deviations are $s_x = 1.73$ and $s_y = 1.00$.

See Chapter 5 for step-by-step calculations.

3. The differences found in Step 3 multiplied together are: $(3 - 4)(2 - 3) = (-1)(-1) = +1$; $(3 - 4)(3 - 3) = (-1)(0) = 0$; $(6 - 4)(4 - 3) = (2)(1) = +2$.

4. Adding the Step 3 results, you get $1 + 0 + 2 = 3$.

5. Dividing by $s_x * s_y$ gives you $3 \div (1.73 * 1.00) = 3 \div 1.73 = 1.73$.

6. Now divide the Step 5 result by $3 - 1$ (which is 2), and you get the correlation $r =$

Interpreting the correlation



The correlation r is always between +1 and -1. To interpret various values of r (no hard and fast rules here, just Rumsey's rule of thumb), see which of the following values your correlation is closest to:

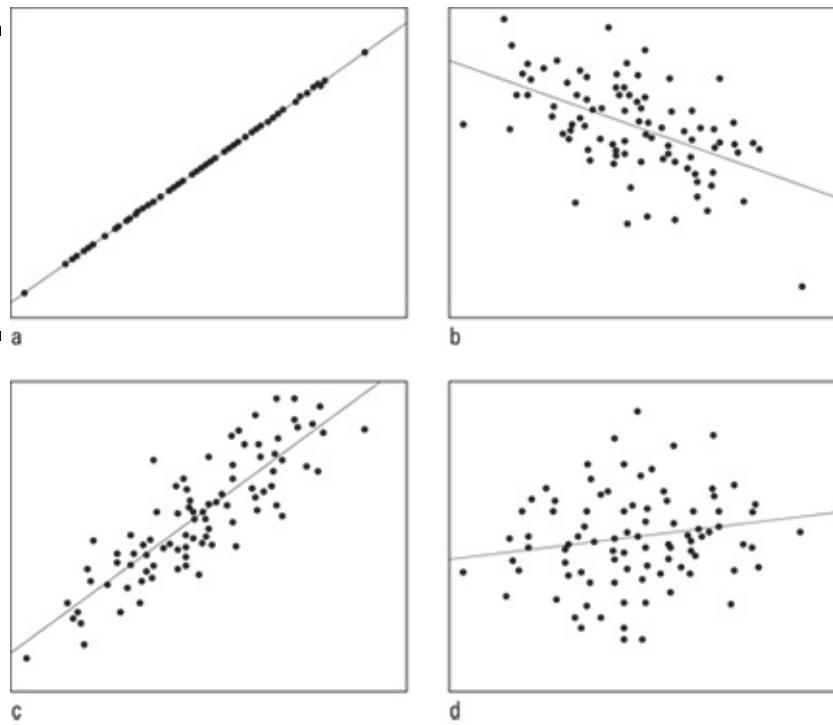
- ✓ **Exactly -1:** A perfect downhill (negative) linear relationship
- ✓ **-0.70:** A strong downhill (negative) linear relationship
- ✓ **-0.50:** A moderate downhill (negative) relationship
- ✓ **-0.30:** A weak downhill (negative) linear relationship
- ✓ **0:** No linear relationship
- ✓ **+0.30:** A weak uphill (positive) linear relationship
- ✓ **+0.50:** A moderate uphill (positive) relationship
- ✓ **+0.70:** A strong uphill (positive) linear relationship
- ✓ **Exactly +1:** A perfect uphill (positive) linear relationship



If the scatterplot doesn't indicate there's at least somewhat of a linear relationship, the correlation doesn't mean much. Why measure the amount of linear relationship if there isn't enough of one to speak of? However you can take the idea of no linear relationship two ways: 1) If no relationship at all exists, calculating the correlation doesn't make sense because correlation only applies to linear relationships; and 2) If a strong relationship exists but it's not linear, the correlation may be misleading, because in some cases a strong curved relationship exists yet the correlation turns out to be strong. That's why it's critical to examine the scatterplot first.

Figure 18-2 shows examples of what various correlations look like, in terms of the strength and direction of the relationship. Figure 18-2a shows a correlation of +1, Figure 18-2b shows a correlation of -0.50, Figure 18-2c shows a correlation of +0.85, and Figure 18-2d shows a correlation of +0.15. Comparing Figures 18-2a and c, you see Figure 18-2a is a perfect uphill straight line, and Figure 18-2c shows a very strong uphill linear pattern. Figure 18-2b is going downhill but the points are somewhat scattered in a wider band, showing a linear relationship is present, but not as strong as in Figures 18-2a and 18-2c. Figure 18-2d doesn't show much of anything happening (and it shouldn't, since its correlation is very close to 0).

Figure 18-2:
Scatterplots
with
correlations
of a) +1.00; b)
-0.50;
c) +0.85; and
d) +0.15.



Many folks make the mistake of thinking that a correlation of -1 is a bad thing, indicating no relationship. Just the opposite is true! A correlation of -1 means the data are lined up in a perfect straight line, the strongest linear relationship you can get. The “ $-$ ” (minus) sign just happens to indicate a negative relationship, a downhill line.



How close is close enough to -1 or $+1$ to indicate a strong enough linear relationship? Most statisticians like to see correlations beyond at least $+0.5$ or -0.5 before getting too excited about them. Don’t expect a correlation to always be 0.99 however; remember, this is real data, and real data aren’t perfect.

For my subset of the cricket chirps versus temperature data from the earlier section “Picturing a Relationship with a Scatterplot,” I calculated a correlation of 0.98 , which is almost unheard of in the real world (these crickets are *good!*).

Examining properties of the correlation



Here are several important properties of the correlation coefficient:

- ✓ The correlation is always between -1 and $+1$, as I explain in the preceding section.
- ✓ The correlation is a unitless measure, which means that if you change the units of X or Y , the correlation won’t change. For example, changing the temperature from Fahrenheit to Celsius won’t affect the correlation between the frequency of chirps (X) and the outside temperature (Y).

- ✓ The variables X and Y can be switched in the data set without changing the correlation. For example, if height and weight have a correlation of 0.53, weight and height have the same correlation.

Working with Linear Regression

In the case of two numerical variables X and Y , when at least a moderate correlation has been established through both the correlation and the scatterplot, you know they have some type of linear relationship. Researchers often use that relationship to predict the (average) value of Y for a given value of X using a straight line. Statisticians call this line the *regression line*. If you know the slope and the y -intercept of that regression line, then you can plug in a value for X and predict the average value for Y . In other words, you predict (the average) Y from X . In the following sections, I provide the basics of understanding and using the linear regression equation (I explain how to make predictions with linear regression later in this chapter).



Never do a regression analysis unless you have already found at least a moderately strong correlation between the two variables. (My rule of thumb is it should be at or beyond either positive or negative 0.50, but other statisticians may have different criteria.) I've seen cases where researchers go ahead and make predictions when a correlation is as low as 0.20! By anyone's standards, that doesn't make sense. If the data don't resemble a line to begin with, you shouldn't try to use a line to fit the data and make predictions (but people still try).

Figuring out which variable is X and which is Y

Before moving forward to find the equation for your regression line, you have to identify which of your two variables is X and which is Y . When doing correlations (as I explain earlier in this chapter), the choice of which variable is X and which is Y doesn't matter, as long as you're consistent for all the data. But when fitting lines and making predictions, the choice of X and Y does make a difference.



So how do you determine which variable is which? In general, Y is the variable that you want to predict, and X is the variable you are using to make that prediction. In the earlier cricket chirps example, you are using the number of chirps to predict the temperature. So in this case the variable Y is the temperature, and the variable X is the number of chirps. Hence Y can be predicted by X using the equation of a line if a strong enough linear relationship exists.



Statisticians call the X -variable (cricket chirps in my earlier example) the *explanatory variable*, because if X changes, the slope tells you (or explains) how much Y is expected to change in response. Therefore, the Y variable is called the *response variable*. Other names for X and Y include the *independent* and *dependent* variables, respectively.

Checking the conditions



In the case of two numerical variables, you can come up with a line that enables you to predict Y from X , if (and only if) the following two conditions from the previous sections are met:

- ✓ The scatterplot must form a linear pattern.
- ✓ The correlation, r , is moderate to strong (typically beyond 0.50 or -0.50).

Some researchers actually don't check these conditions before making predictions. Their claims are not valid unless the two conditions are met.

But suppose the correlation is high; do you still need to look at the scatterplot? Yes. In some situations the data have a somewhat curved shape, yet the correlation is still strong; in these cases making predictions using a straight line is still invalid. Predictions need to be made based on a curve. (This topic is outside the scope of this book; if you are interested, see *Statistics II For Dummies*, where I tackle nonlinear relationships.)

Calculating the regression line

For the crickets and temperature data, you can see that the scatterplot in Figure 18-1 shows a linear pattern. The correlation between cricket chirps and temperature was found earlier in this chapter to be very strong ($r = 0.98$). You now can find one line that best fits the data (in terms of having the smallest overall distance to the points). Statisticians call this technique for finding the best-fitting line a *simple linear regression analysis using the least squares method*.



The formula for the *best-fitting line* (or *regression line*) is $y = mx + b$, where m is the slope of the line and b is the y -intercept. This equation itself is the same one used to find a line in algebra; but remember, in statistics the points don't lie perfectly on a line — the line is a model around which the data lie if a strong linear pattern exists.

- ✓ The *slope* of a line is the change in Y over the change in X . For example, a slope of

10/3 means as the x -value increases (moves right) by 3 units, the y -value moves up by 10 units on average.

- ✓ The *y-intercept* is that place on the y -axis where the value of x is zero. For example, in the equation $2x - 6$, the line crosses the y -axis at the point -6 . The coordinates of this point are $(0, -6)$; when a line crosses the y -axis, the x -value is always 0.



To come up with the best-fitting line, you need to find values for m and b that fit the pattern of data the best, for your given criteria. Different criteria exist and can lead to other lines, but the criteria I use in this book (and in all introductory level statistics courses in general) is to find the line that minimizes what statisticians call the *sum of squares for error (SSE)*. The SSE is the sum of all the squared differences from the points on the proposed line to the actual points in the data set. The line with the lowest possible SSE wins and its equation is used as the best-fitting line. This process is where the name *the least-squares method* comes from.

You may be thinking that you have to try lots and lots of different lines to see which one fits best. Fortunately, you have a more straightforward option (although eyeballing a line on the scatterplot does help you think about what you'd expect the answer to be). The best-fitting line has a distinct slope and y -intercept that can be calculated using formulas (and, I may add, these formulas aren't too hard to calculate).



To save a great deal of time calculating the best fitting line, first find the “big five,” five summary statistics that you’ll need in your calculations:

1. The mean of the x values (denoted \bar{x})
2. The mean of the y values (denoted \bar{y})
3. The standard deviation of the x values (denoted s_x)
4. The standard deviation of the y values (denoted s_y)
5. The correlation between X and Y (denoted r)

Finding the slope

The formula for the slope, m , of the best-fitting line is

$$m = r \left(\frac{s_y}{s_x} \right)$$

where r is the correlation between X and Y , and s_x and s_y are the standard deviations of the x -values and the y -values, respectively. You simply divide s_y by s_x and multiply the result by r .

Note that the slope of the best-fitting line can be a negative number because the

correlation can be a negative number. A negative slope indicates that the line is going downhill. For example, an increase in police officers is related to a decrease in the number of crimes in a linear fashion; the correlation and hence the slope of the best-fitting line is negative in this case.



The correlation and the slope of the best-fitting line are not the same. The formula for slope takes the correlation (a unitless measurement) and attaches units to it. Think of $s_y \div s_x$ as the variation (resembling change) in Y over the variation in X , in units of X and Y . For example, variation in temperature (degrees Fahrenheit) over the variation in number of cricket chirps (in 15 seconds).

Finding the y-intercept

The formula for the y -intercept, b , of the best-fitting line is $b = \bar{y} - m\bar{x}$, where \bar{x} and \bar{y} are the means of the x -values and the y -values, respectively, and m is the slope (the formula for which is given in the preceding section).



So to calculate the y -intercept, b , of the best-fitting line, you start by finding the slope, m , of the best-fitting line using the steps listed in the preceding section. You then multiply m by \bar{x} and subtract your result from \bar{y} .



Always calculate the slope before the y -intercept. The formula for the y -intercept contains the slope!

Interpreting the regression line

Even more important than being able to calculate the slope and y -intercept to form the best-fitting regression line is the ability to interpret their values; I explain how to do so in the following sections.

Interpreting the slope

The slope is interpreted in algebra as *rise over run*. If, for example, the slope is 2, you can write this as $2/1$ and say that as you move from point to point on the line, as the value of the X variable increases by 1, the value of the Y variable increases by 2. In a regression context, the slope is the heart and soul of the equation because it tells you how much you can expect Y to change as X increases.

In general, the units for slope are the units of the Y variable per units of the X variable. It's a ratio of change in Y per change in X . Suppose in studying the effect of dosage level

in milligrams (mg) on systolic blood pressure (mmHg), a researcher finds that the slope of the regression line is -2.5 . You can write this as **$-2.5/1$** and say that systolic blood pressure is expected to decrease by 2.5 mmHg on average per 1 mg increase in drug dosage.



Always make sure to use proper units when interpreting slope. If you don't consider units, you won't really see the connection between the two variables at hand. For example if Y is exam score and X = study time, and you find the slope of the equation is 5, what does this mean? Not much without any units to draw from. Including the units, you see you get an increase of 5 points (change in Y) for every 1 hour increase in studying (change in X). Also be sure to watch for variables that have more than one common unit, such as temperature being in either Fahrenheit or Celsius; know which unit is being used.

If using a 1 in the denominator of slope is not super-meaningful to you, you can multiply the top and bottom by any number (as long as it's the same number) and interpret it that way instead. In the systolic blood pressure example, instead of writing slope as **$-2.5/1$** and interpreting it as a drop of 2.5 mmHg per 1 mg increase of the drug, you can multiply the top and bottom by 10 to get **$-25/10$** and say an increase in dosage of 10 mg results in a decrease in systolic blood pressure of 25 mmHg.

Interpreting the y-intercept

The y -intercept is the place where the regression line $y = mx + b$ crosses the y -axis where $x = 0$, and is denoted by b (see the earlier section “Finding the y -intercept”). Sometimes the y -intercept can be interpreted in a meaningful way, and sometimes not. This uncertainty differs from slope, which is always interpretable. In fact, between the two elements of slope and y -intercept, the slope is the star of the show, with the y -intercept serving as the less-famous but still noticeable sidekick.



At times the y -intercept makes no sense. For example, suppose you use rain to predict bushels per acre of corn. You know if the data set contains a point where rain is 0, the bushels per acre must be 0 as well. As a result, if the regression line crosses the y -axis somewhere else besides 0 (and there is no guarantee it will cross at 0 — it depends on the data), the y -intercept will make no sense. Similarly, in this context a negative value of y (corn production) cannot be interpreted.

Another situation where you can't interpret the y -intercept is when data are not present near the point where $x = 0$. For example, suppose you want to use students' scores on Midterm 1 to predict their scores on Midterm 2. The y -intercept represents a prediction for Midterm 2 when the score on Midterm 1 is 0. You don't expect scores on a midterm to be at or near 0 unless someone didn't take the exam, in which case her score wouldn't be included in the first place.

Many times, however, the y -intercept is of interest to you, it has meaning, and you have data collected in the area where $x = 0$. For example, if you're predicting coffee sales at football games in Green Bay, Wisconsin, using temperature, some games get cold enough to have temperatures at or even below 0 degrees Fahrenheit, so predicting coffee sales at these temperatures makes sense. (As you may guess, they sell more and more coffee as the temperature dips.)

Putting it all together with an example: The regression line for the crickets

In the earlier section "Picturing a Relationship with a Scatterplot," I introduce the example of cricket chirps related to temperature. The "big five" statistics, which I explain in "Calculating the regression line," are shown in Table 18-2 for the subset of cricket data. (**Note:** I'm rounding for ease of explanation only.)

Table 18-2 "Big Five" Statistics for the Cricket Data			
Variable	Mean	Standard Deviation	Correlation
Number of chirps (x)	$\bar{x} = 26.5$	$s_x = 7.4$	$r = +0.98$
Temp (y)	$\bar{y} = 67$	$s_y = 6.8$	

The slope, m , for the best-fitting line for the subset of cricket chirps versus temperature data is $m = r \frac{s_y}{s_x} = 0.98 \left(\frac{6.8}{7.4} \right) = 0.90$. So as the number of chirps increases by 1 chirp per 15 seconds, the temperature is expected to increase by 0.90 degrees Fahrenheit on average. To get a more meaningful interpretation, you can multiply the top and bottom of the slope by 10 and say as chirps increase by 10 (per 15 seconds) temperature increases 9 degrees Fahrenheit.

Now, to find the y -intercept, b , you take $\bar{y} - m\bar{x}$, or $67 - (0.90)(26.5) = 43.15$. So the best-fitting line for predicting temperature from cricket chirps based on the data is $y = 0.90x + 43.15$, or temperature (in degrees Fahrenheit) = $0.90 * (\text{number of chirps in 15 seconds}) + 43.2$. Now can you use the y -intercept to predict temperature when no chirping is going on at all? Because no data was collected at or near this point, you cannot make predictions for temperature in this area. You can't predict temperature using crickets if the crickets are silent.

Making Proper Predictions

After you have determined a strong linear relationship and you find the equation of the best fitting line using $y = mx + b$, you use that line to predict (the average) y for a given x -value. To make predictions, you plug the x -value into the equation and solve for y . For

example, if your equation is $y = 2x + 1$ and you want to predict y for $x = 1$, then plug 1 into the equation for x to get $y = 2(1) + 1 = 3$.

Keep in mind that you choose the values of X (the explanatory variable) that you plug in; what you predict is Y , the response variable, which totally depends on X . By doing this, you are using one variable that you can easily collect data on to predict a Y variable that is difficult or not possible to measure. This process works well as long as X and Y are correlated. This concept is the big idea of regression.

Using the examples from the previous section, the best-fitting line for the crickets is $y = 0.90x + 43.2$. Say you're camping outside, listening to the crickets, and remember you can predict temperature by counting cricket chirps. You count 35 chirps in 15 seconds, put in 35 for x , and find that $y = 0.9(35) + 43.2 = 74.7$. (Yeah, you memorized the formula before you went camping just in case you needed it.) So because the crickets chirped 35 times in 15 seconds, you figure the temperature is probably about 75 degrees Fahrenheit.



Just because you have a regression line doesn't mean you can plug in *any* value for X and do a good job of predicting Y . Making predictions using x -values that fall outside the range of your data is a no-no. Statisticians call this *extrapolation*; watch for researchers who try to make claims beyond the range of their data.

For example, in the chirping data, no data is collected for fewer than 18 chirps or more than 39 chirps per 15 seconds (refer to Table 18-1). If you try to make predictions outside this range, you are going into uncharted territory; the farther outside this range you go with your x -values, the more dubious your predictions for y will get. Who's to say the line still works outside of the area where data were collected? Do you really think that crickets will chirp faster and faster without limit? At some point they would either pass out or burn up! And what does a negative number of chirps really mean? (Is this similar to asking what the sound of one hand clapping is?)



Be aware that not every data point will necessarily fit the regression line well, even if the correlation is high. A point or two may fall outside the overall pattern of the rest of the data; such points are called *outliers*. One or two outliers probably won't affect the overall fit of the regression line much, but in the end you can see that the line didn't do well at those specific points.

The numerical difference between the predicted value of y from the line and the actual y -value you got from your data is called a *residual*. Outliers have large residuals compared to the rest of the points; they are worth investigating to see if there was an error in the data at those points or if there is something particularly interesting in the data to follow up on. (I give a much more detailed look at residuals in the book *Statistics II For Dummies*.)

Explaining the Relationship: Correlation versus Cause and Effect

Scatterplots and correlations identify and quantify relationships between two variables. However, if a scatterplot shows a definite pattern and the data are found to have a strong correlation, that doesn't necessarily mean that a cause-and-effect relationship exists between the two variables. A *cause-and-effect relationship* is one where a change in one variable (in this case X) causes a change in another variable (in this case Y). (In other words, the change in Y is not only associated with a change in X , but also directly caused by X .)

For example, suppose a well-controlled medical experiment is conducted to determine the effects of dosage of a certain drug on blood pressure. (See a total breakdown of experiments in Chapter 17.) The researchers look at their scatterplot and see a definite downhill linear pattern; they calculate the correlation, and it's strong. They conclude that increasing the dosage of this drug causes a decrease in blood pressure. This cause-and-effect conclusion is okay because they controlled for other variables that could affect blood pressure in their experiment, such as other drugs taken, age, general health, and so on.

However, if you made a scatterplot and examined the correlation between ice cream consumption versus murder rates in New York City, you would also see a strong linear relationship (this one is uphill). Yet no one would claim that more ice cream consumption causes more murders to occur.

What's going on here? In the first case, the data were collected through a well-controlled medical experiment, which minimizes the influence of other factors that may affect blood pressure. In the second example, the data were based just on observation, and no other factors were examined. Researchers subsequently found out that this strong relationship exists because increases in murder rates and ice cream sales are both related to increases in temperature. Temperature in this case is called a *confounding variable*; it affects both X and Y but was not included in the study (see Chapter 17).



Whether two variables are found to be causally associated depends on how the study was conducted. I've seen many instances in which people try to claim cause-and-effect relationships just by looking at scatterplots or correlations. Why would they do this? Because they want to believe it (in other words for them it's "believing is seeing," rather than the other way around). Beware of this tactic. In order to establish cause and effect, you need to have a well-designed experiment or a boatload of observational studies. If someone is trying to establish a cause-and-effect relationship by showing a chart or graph, dig deeper to find out how the study was designed and how the data were collected, and evaluate the study appropriately.

using the criteria outlined in Chapter 17.

The need for a well-designed experiment in order to claim cause and effect is often ignored by some researchers and members of the media, who give us headlines such as “Doctors can lower malpractice lawsuits by spending more time with patients.” In reality, it was found that doctors who have fewer lawsuits are the type who spend a lot of time with patients. But that doesn’t mean taking a bad doctor and having him spend more time with his patients will reduce his malpractice suits; in fact, spending more time with them may create even more problems.

Chapter 19

Two-Way Tables and Independence

In This Chapter

- ▶ Setting up two-way tables with categorical variables
 - ▶ Delving into marginal, joint, and conditional distributions
 - ▶ Checking for independence and dependence
 - ▶ Having perspective on the results of two-way tables
-

Categorical variables place individuals into groups based on certain characteristics, behaviors, or outcomes, such as whether you ate breakfast this morning (yes, no) or political affiliation (Democrat, Republican, Independent, “other”). Oftentimes people look for relationships between two categorical variables; hardly a day goes by that you don’t hear about another relationship that’s reported to have been found.

Here are just a few examples I found on the Internet recently:

- ✓ Dog owners are more likely to take their animal to the vet than cat owners.
- ✓ Heavy use of social-networking Web sites in teens is linked to depression.
- ✓ Children who play more video games do better in science classes.

With all this information being given to you about variables that are related, how do you decide what to believe? For example, does heavy use of social-networking Web sites cause depression, or is it the other way around? Or perhaps a third variable out there is related to both of them, such as problems in the home.

In this chapter, you see how to organize and analyze data from two categorical variables. You find out how to use proportions to make comparisons and look at overall patterns and how to check for independence of two categorical variables. You see how to describe dependent relationships appropriately and to evaluate results claiming to indicate cause-and-effect relationships, making predictions, and/or projecting their results to a population.

Organizing a Two-Way Table

To explore links between two categorical variables, you first need to organize the data that’s been collected, and a table is a great way to do that. A *two-way table* classifies

individuals into groups based on the outcomes of two categorical variables (for example, gender and opinion).

Suppose your local community developers are building a campground, and they've decided pets will be allowed as long as they're on a leash. They are now trying to decide whether the campground should have a separate section for pets. You have a hunch that non-pet campers in the area may be more in favor of a separate pet area than pet campers, so you decide to find out what the members of the camping community think. You randomly select 100 campers from the local area and conduct a pet camping survey, recording each person's opinion on having a pet section (yes, no) and if they camp with pets (yes, no). You now have a spreadsheet with 100 rows of data, one for each person you surveyed. Each row has two pieces of data: one column for whether the person is a pet camper (yes, no) and one column for that person's opinion on having a pet section (support, oppose). Suppose the first 10 rows of your data set look like what's shown in Table 19-1.

Table 19-1 First 10 Rows of Data from the Pet Camping Survey

<i>Person</i>	<i>Pet Camper?</i>	<i>Opinion on a Separate Pet Section</i>
1	Yes	Oppose
2	Yes	Oppose
3	Yes	Support
4	No	Support
5	No	Support
6	Yes	Support
7	No	Oppose
8	No	Support
9	Yes	Support
10	No	Oppose

From this small portion of your data set, you can start to break it down yourself. For example, looking at column 2 results, you see that half the respondents ($5 \div 10 = 0.50$) camp with pets and the other half do not. Of those who camp with pets (that is, of those five people who have a yes in column 2), three of them (60%) support having a separate section; and the same results are true for non-pet campers. These results from these 10 campers likely don't apply to all 100 campers surveyed; however, if you tried to examine the raw data from all 100 rows of this data set by hand, you wouldn't make much progress in seeing patterns without a lot of hard work.

In order to get a handle on what's happening in a large data set when you are examining two categorical variables, you organize your data into a two-way table. The following sections take you through it.

Setting up the cells



A two-way table organizes categorical data from two variables by using rows to represent one variable (such as pet camping — yes or no) and columns to represent the other variable (such as opinion on a pet section — support or oppose). Each person appears exactly once in the table.

Continuing with the camping example I start earlier in this chapter, in Table 19-2 I summarize the results from all 100 campers surveyed.

**Table 19-2 Two-Way Table of Pet Camping Survey Data
(All 100 Rows)**

	<i>Support Separate Pet Section</i>	<i>Oppose Separate Pet Section</i>
<i>Pet Camper</i>	20	10
<i>Non-Pet Camper</i>	55	15

Table 19-2 has $2 * 2 = 4$ numbers in it. These numbers represent the *cells* of the two-way table; each one represents an intersection of a row and column. The cell in the upper left corner of the table represents the 20 people who are pet campers supporting a pet section. In the upper right cell 10 people are pet campers opposing a pet section. In the lower left are the 55 non-pet campers who want a pet section; the 15 people in the lower right are non-pet campers opposing a pet section.

Figuring the totals



Before getting to the nitty-gritty analysis of a two-way table in the later section “Interpreting Results from a Two-Way Table,” you calculate some totals and add them to the table for later reference. You summarize each variable separately by calculating the *marginal totals*, which represent the total number in each row (for the first variable) and the total number in each column (for the second variable). The *marginal row totals* form an additional column on the right side of the table, and the *marginal column totals* form an additional row on the bottom of the table.

For example, in Table 19-2 in the preceding section, the marginal row total for row 1, the number of pet campers, is $20 + 10 = 30$; the marginal row total for non-pet campers (row 2) is $55 + 15 = 70$. The marginal column total for those wanting a pet section (column 1) is $20 + 55 = 75$; and the marginal column total for those not wanting a separate section (column 2) is $10 + 15 = 25$.



The *grand total* is the total of all the cells in the table and is equal to the sample size. (Note the marginal totals are not included in the grand total, only the cells.) The grand total sits in the lower right-hand corner of the two-way table. In this example, the grand total is $20 + 10 + 55 + 15 = 100$. Table 19-3 shows the marginal row and column totals and the grand total for the pet camping survey data.

The marginal row totals always sum to the grand total, because everyone in the survey either camps with a pet or they don't. In the last column of Table 19-3 you see that $30 + 70 = 100$. Similarly the marginal column totals always sum to the grand total; everyone in the survey either wants a pet section or they don't; in the last row of Table 19-3 you see $75 + 25 = 100$.

Table 19-3 Two-Way Table of Pet Camping Survey Data, Including Marginal Totals

	<i>Support Separate Pet Section</i>	<i>Oppose Separate Pet Section</i>	<i>Marginal Row Totals</i>
Pet Camper	20	10	$20 + 10 = 30$
Non-Pet Camper	55	15	$55 + 15 = 70$
Marginal Column Totals	$20 + 55 = 75$	$10 + 15 = 25$	Grand total = 100 ($20 + 10 + 55 + 15$)



When organizing a two-way table, always include the marginal totals and the grand total. It gets you off on the right foot when analyzing the data.

Interpreting Two-Way Tables

After the two-way table is set up (with the help of the information in the previous section), you calculate percents to explore the data to answer your research questions. Here are some questions of interest from the camping data earlier in this chapter (each question will be handled in the following sections, respectively):

- ✓ What percentage of the campers are in favor of a pet section?
- ✓ What percentage of the campers are pet campers who support a pet section?
- ✓ Do more non-pet campers support a pet section, compared to pet campers?

The answers to these (and any other) questions about the data come from finding and working with the proportions, or percentages, of individuals within certain parts of the table. This process involves calculating and examining what statisticians call

distributions. A distribution in the case of a two-way table is a list of all the possible outcomes for one variable or a combination of variables, along with their corresponding proportions (or percentages).

For example, the distribution for the pet camping variable lists the percentages of people who do and do not camp with pets. The distribution for the combination of the pet camping variable (yes, no) and the opinion variable (support, oppose) lists the percentages of: 1) pet campers who support a pet section; 2) pet campers who oppose a pet section; 3) non-pet campers who support a pet section; and 4) the non-pet campers who oppose a pet section.



For any distribution, all the percentages must sum to 100%. If you're using proportions (decimals), they must sum to 1.00. Each individual has to be somewhere, and he can't be in more than one place at one time.

In the following sections, you see how to find three types of distributions, each one helping you to answer its corresponding question in the preceding list.

Singling out variables with marginal distributions

If you want to examine one variable at a time in a two-way table, you don't look in the cells of the table, but rather in the margins. As seen in the earlier section "Figuring the totals," the marginal totals represent the total number in each row (or column) separately. In the two-way table for the pet camping survey (refer to Table 19-3), you see the marginal totals for the pet camping variable (yes/no) in the right-hand column, and you find the marginal totals for the opinion variable (support/oppose) in the bottom row.

If you want to make comparisons between two groups (for example, pet campers versus non-pet campers), however, the results are easier to interpret if you use proportions instead of totals. If 350 people were surveyed, visualizing a comparison is easier if you're told that 60% are in Group A and 40% are in Group B, rather than saying 210 people are in Group A and 140 are in Group B.

To examine the results of a two-way table based on a single variable, you find what statisticians call the *marginal distribution* for that variable. In the following sections, I show you how to calculate and graph marginal distributions.

Calculating marginal distributions



To find a marginal distribution for one variable in a two-way table, you take the marginal total for each row (or column) divided by the grand total.

- ✓ If your variable is represented by the rows (for example, the pet camping variable in Table 19-3), use the marginal row totals in the numerators and the grand total in the denominators. Table 19-4 shows the marginal distribution for the pet camping variable (yes, no).
- ✓ If your variable is represented by the columns (for example, opinion on the pet section policy, shown in Table 19-3), use the marginal column totals for the numerators and the grand total for the denominators. Table 19-5 shows the marginal distribution for the opinion variable (support, oppose).



In either case, the sum of the proportions for any marginal distribution must be 1 (subject to rounding). All results in a two-way table are subject to rounding error; to reduce rounding error, keep at least 2 digits after the decimal point throughout.

Table 19-4 Marginal Distribution for Pet Camping Variable

<i>Pet Camping</i>	<i>Proportion</i>
Yes	$30 \div 100 = 0.30$
No	$70 \div 100 = 0.70$
Total	1.00

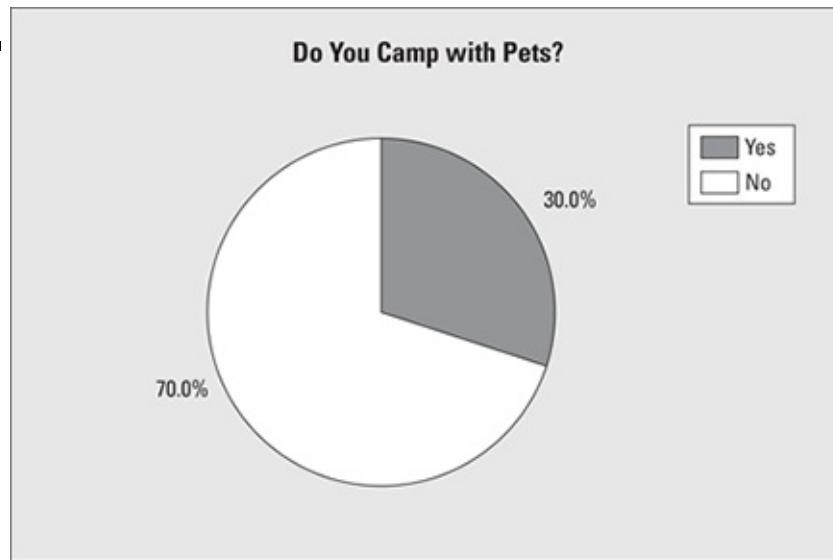
Table 19-5 Marginal Distribution for the Opinion Variable

<i>Opinion</i>	<i>Proportion</i>
Support pet section	$75 \div 100 = 0.75$
Oppose pet section	$25 \div 100 = 0.25$
Total	1.00

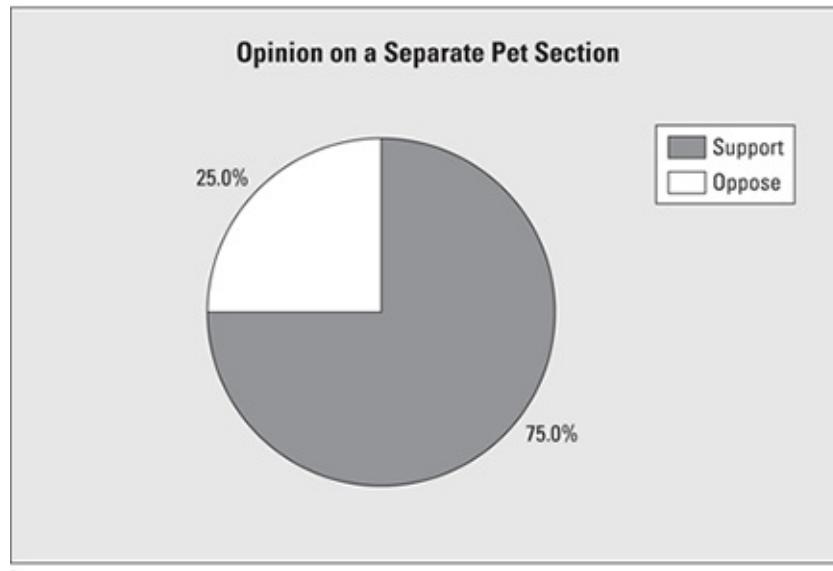
Graphing marginal distributions

You graph a marginal distribution using either a pie chart or a bar graph. Each graph shows the proportion of individuals within each group for a single variable. Figure 19-1a is a pie chart summarizing the pet camping variable, and Figure 19-1b is a pie chart showing the breakdown of the opinion variable. You see that the results of these two pie charts correspond with the marginal distributions in Tables 19-4 and 19-5, respectively.

Figure 19-1:
Pie charts
showing
marginal
distributions
for a) pet
camping
variable; and
b) opinion



a



b

From the results of the two separate marginal distributions for the pet camping and opinion variables, you say that the majority of all the campers in this sample are non-pet campers (70%) and the majority of all the campers in this sample (75%) support the idea of having a pet section.



While marginal distributions show us how each variable breaks down on its own, they don't tell us about the connection between two variables. For the camping example, you know what percentage of all campers support a new pet section, but you can't distinguish the opinions of the pet campers from the non-pet campers. Distributions for making such comparisons are found in the later section, "Comparing groups with conditional distributions."

Examining all groups — a joint distribution

Story time: A certain auto manufacturer conducted a survey to see what characteristics

customers prefer in their small pickup trucks. They found that the most popular color for these trucks was red and the most popular option was four-wheel drive. In response to these results, the company started making more of their small pickup trucks red with four-wheel drive.

Guess what? They struck out; people weren't buying those trucks. Turns out that the customers who bought the red trucks were more likely to be women, and women didn't use four-wheel drive as often as men did. Customers who bought the four-wheel drive trucks were more likely to be men, and they tended to prefer black ones over red ones. So the most popular outcome of the first variable (color) paired with the most popular outcome of the second variable (options on the vehicle) doesn't necessarily add up to the most popular combination of the two variables.



To figure out which combination of two categorical variables contains the highest proportion, you need to compare the cell proportions (for example, the color and vehicle options together) rather than the marginal proportions (the color and vehicle option separately). The *joint distribution* of both variables in a two-way table is a listing of all possible row and column combinations and the proportion of individuals within each group. You use it to answer questions involving two characteristics; such as "What proportion of the voters are Democrat and female?" or, "What percentage of the campers are pet campers who support a pet section?" In the following sections, I show you how to calculate and graph joint distributions.

Calculating joint distributions

A joint distribution shows the proportion of the data that lies in each cell of the two-way table. For the pet camping example, the four row-column combinations are:

- ✓ All campers who camp with pets and support a pet section.
- ✓ All campers who camp with pets and oppose a pet section.
- ✓ All campers who don't camp with pets and support a pet section.
- ✓ All campers who don't camp with pets and oppose a pet section.



The key phrase in all of the proportions mentioned in the preceding list is *all campers*. You are taking the entire group of all campers in the survey and breaking them into four separate groups. When you see the word *all*, think joint distribution. Table 19-6 shows the joint distribution for all campers in the pet camping survey.

Table 19-6 Joint Distribution for the Pet Camping Survey Data

	<i>Support Separate Pet Section</i>	<i>Oppose Separate Pet Section</i>
Camp with Pets	$20 \div 100 = 0.20$	$10 \div 100 = 0.10$
Don't Camp with Pets	$55 \div 100 = 0.55$	$15 \div 100 = 0.15$



To find a joint distribution for a two-way table, you take the cell count (the number of individuals in a cell) divided by the grand total, for each cell in the table. The total of all these proportions should be 1 (subject to rounding error).

To get the numbers in the cells of Table 19-6, take the cells of Table 19-3 and divide by their corresponding grand total (100, in this case). Using the results listed in Table 19-6, you report the following:

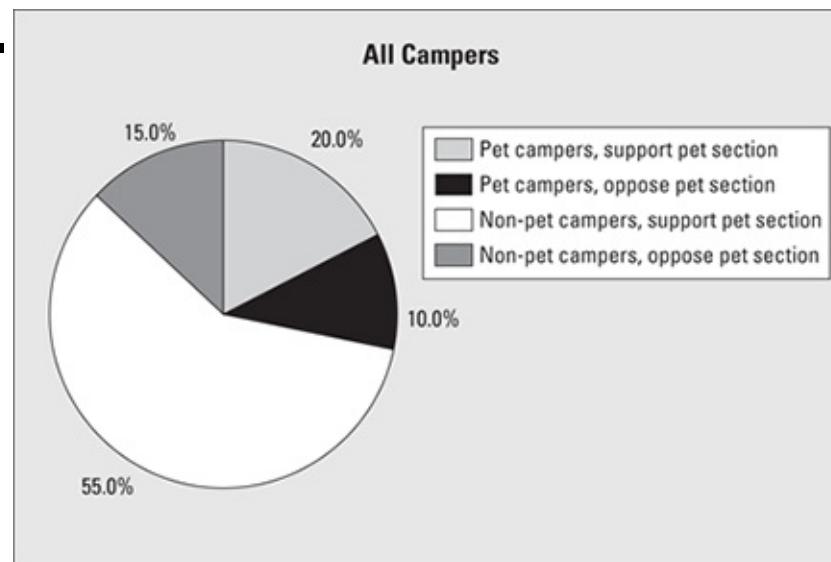
- ✓ 20% of all campers surveyed camp with pets and support a pet section. (See the upper left-hand cell of the table.)
- ✓ 10% of all campers surveyed camp with pets and oppose a pet section. (See the upper right-hand cell of the table.)
- ✓ 55% of all campers surveyed don't camp with pets and do support the pet section policy. (See the lower left-hand cell of the table.)
- ✓ 15% of all campers surveyed don't camp with pets and oppose the pet section policy. (See the lower right-hand cell of the table.)

Adding all the proportions shown in Table 19-6, you get $0.20 + 0.10 + 0.55 + 0.15 = 1.00$. Every camper shows up in one and only one of the cells of the table.

Graphing joint distributions

To graph a joint distribution from a two-way table, you make a single pie chart with four slices, representing each proportion of the data that falls within a row-column combination. Groups containing more individuals get a bigger piece of the overall pie, and hence get more weight when all the votes are counted up. Figure 19-2 is a pie chart showing the joint distribution for the pet camping survey data.

Figure 19-2:
Pie chart
showing the
joint
distribution of
the pet
camping and
opinion



From the pie chart shown in Figure 19-2, you see some results that stand out. The majority of campers in this sample (0.55 or 55%) don't camp with pets and support a separate section for pets. The smallest slice of the pie represents those campers who camp with pets and are opposed to a separate section for pets (0.10 or 10%).

A joint distribution gives you a breakdown of the entire group by both variables at once and allows you to compare the cells to each other and to the whole group. The results in Figure 19-2 show that if they were asked to vote today as to whether or not to have a pet section, when all the votes were added up, most of the weight would be placed on the opinions of non-pet campers, because they make up the majority of campers in the survey (70%, according to Table 19-4), and the pet campers would have less of a voice, because they are a smaller group (30%).



A limitation of a joint distribution is that you can't fairly compare two groups to each other (for example pet campers versus non-pet campers) because the joint distribution puts more weight on larger groups. The next section shows how to fairly compare the groups in a two-way table.

Comparing groups with conditional distributions

You need a different type of distribution other than a joint distribution to compare the results from two groups (for example comparing opinions of pet campers versus non-pet campers). *Conditional distributions* are used when looking for relationships between two categorical variables; the individuals are first split into the groups you want to compare (for example, pet campers and non-pet campers); then the groups are compared based on their opinion on a pet section (yes, no). In the following sections, I explain how to calculate and graph conditional distributions.

Calculating conditional distributions



To find conditional distributions for the purpose of comparison, first split the individuals into groups according to the variable you want to compare. Then for each group, take the cell count (the number of individuals in a particular cell) divided by the marginal total for that group. Do this for all the cells in that group. Now repeat for the other group, using its marginal total as the denominator and the cells within its group as the numerators. (See the earlier section “Figuring the totals” for more about marginal totals.) You now have two conditional distributions, one for each group, and you fairly compare the results for the two groups.

For the pet camping survey data example (earlier in this chapter), you compare the opinions of two groups: pet campers and non–pet campers; in statistical terms you want to find the conditional distributions of opinion based on the pet camping variable. That means you split the individuals into the pet camper and non–pet camper groups, and then for each group, you find the percentages of who supports and opposes the new pet section. Table 19-7 shows these two conditional distributions in table form (working off Table 19-3).

Table 19-7
Conditional Distributions of Opinion for Pet Campers versus Non–Pet Campers

	<i>Support Pet Section Policy</i>	<i>Oppose Pet Section Policy</i>	<i>Total</i>
Pet Campers	$20 \div 30 = 0.67$	$10 \div 30 = 0.33$	1.00
Non–Pet Campers	$55 \div 70 = 0.79$	$15 \div 70 = 0.21$	1.00



Notice that Table 19-7 differs from Table 19-6 in the earlier section “Calculating joint distributions” in terms of how the values in the table add up. This represents the key difference between a joint distribution and a conditional distribution that allows you to make fair comparisons using the conditional distribution:

- ✓ In Table 19-6, the proportions in the cells of the entire table sum to 1 because the entire group is broken down by both variables at once in a joint distribution.
- ✓ In Table 19-7, the proportions in each row of the table sum to 1 because each group is treated separately in a conditional distribution.

Graphing conditional distributions

One effective way to graph conditional distributions is to make a pie chart for each group (for example, one for pet campers and one for non–pet campers) where each pie chart shows the results of the variable being studied (opinion: yes or no).

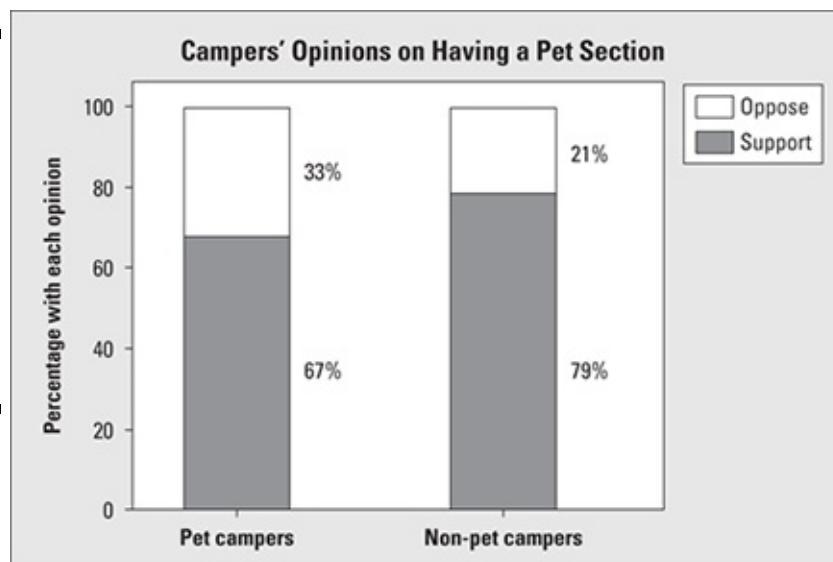
Another method is to use a stacked bar graph. A *stacked bar graph* is a special bar graph

where each bar has a height of 1 and represents an entire group (one bar for pet campers and one bar for non-pet campers). Each bar shows how that group breaks down regarding the other variable being studied (opinion: yes or no).

Figure 19-3 is a stacked bar graph showing two conditional distributions. The first bar is the conditional distribution of opinion for the pet camping group (row 1 of Table 19-7) and the second bar represents the conditional distribution of opinion for the non-pet camping group (row 2 of Table 19-7).

Using Table 19-7 and Figure 19-3, first look at the opinions of each group. More than 50% of the pet campers support the pet section (the exact number rounds to 67%), so you say the majority of pet campers support a pet section. Similarly, the majority of non-pet campers (about 79%, way more than half) support a pet section.

Figure 19-3:
Stacked bar graph showing the conditional distributions of opinion for pet campers and non-pet campers.



Now you compare the opinions of the two groups by comparing the percentage of supporters in the pet camping group (67%) to the percentage of supporters in the non-pet camping group (79%). While both groups have a majority of supporters of the pet section, you see more of the non-pet campers support the policy than pet campers ($79\% > 67\%$). By comparing the conditional distributions, you've found that a relationship appears to exist between opinion and pet camping, and your original hunch that non-pet campers in the area may be more in favor of a separate pet area than pet campers is correct, based on this data.



The difference in the results found in Figure 19-3 isn't as large as you may have thought by looking at the joint distribution in Figure 19-2. The conditional distribution takes into account and adjusts for the number in each group being compared, while the joint distribution puts everyone in the same boat. That's why you need conditional distributions to make fair comparisons.



When making my conclusions regarding the pet-camping data, the operative words I use are “a relationship *appears* to exist.” The results of the pet camping survey are based on only your sample of 100 campers. To be able to generalize these results to the whole population of pet campers and non-pet campers in this community (which is really what you want to do), you need to take into account that these sample results will vary, and when they do vary, will they still show the same kind of difference? That’s what a hypothesis test will tell you (all the details are in Chapter 14).



To conduct a hypothesis test for a relationship between two categorical variables (when each variable has only two categories, like yes/no or male/female), you either do a test for two proportions (see Chapter 15) or a Chi-square test (which is covered in my book *Statistics II For Dummies*, also published by Wiley). If one or more of your variables have more than two categories, such as Democrats/Republicans/Other, you must use the Chi-square test to test for independence in the population.



Be mindful that you may run across a report in which someone is trying to give the appearance of a stronger relationship than really exists, or trying to make a relationship less obvious by how the graphs are made. With pie charts, the sample size often is not reported, leading you to believe the results are based on a large sample when they may not be. With bar graphs, they stretch or shrink the scale to make differences appear larger or smaller, respectively. (See Chapter 6 for more information on misleading graphs of categorical data.)

Checking Independence and Describing Dependence

The main reason researchers collect data on two categorical variables is to explore possible relationships or connections between the variables. For example, if a survey finds that more females than males voted for the incumbent president in the last election, then you conclude that gender and voting outcome are related. If a relationship between two categorical variables has been found (that is, the results from the two groups are different), then statisticians say they’re *dependent*.

However, if you find that the percentage of females who voted for the incumbent is the same as the percentage of males who voted for the incumbent, then the two variables (gender and voting for the incumbent) have no relationship and statisticians say those two variables are *independent*. In this section, you find out how to check for independence and describe relationships found to be dependent.

Checking for independence

Two categorical variables are *independent* if the percentages for the second variable (typically representing the results you want to compare, such as support or oppose) do not differ based on the first variable (typically representing the groups you want to compare, such as men versus women). You can check for independence with the methods that I cover in this section.

Comparing the results of two conditional distributions



Two categorical variables are *independent* if the conditional distributions are the same for all groups being compared. The variables are independent because breaking them down and comparing them by group doesn't change the results. In the election example I introduced at the beginning of "Checking Independence and Describing Dependence," independence means the conditional distribution for opinion is the same for the males and the females.

Suppose you do a survey of 200 voters to see if gender is related to whether they voted for the incumbent president, and you summarize your results in Table 19-8.

Table 19-8 Results of Election Survey

	<i>Voted for Incumbent President</i>	<i>Didn't Vote for Incumbent President</i>	<i>Marginal Row Totals</i>
<i>Males</i>	44	66	110
<i>Females</i>	36	54	90
<i>Marginal Column Totals</i>	80	120	Grand total = 200

To see whether gender and voting are independent, you find the conditional distribution of voting pattern for the males and the conditional distribution of voting pattern for the females. If they're the same, you've got independence; if not, you've got dependence. These two conditional distributions have been calculated and appear in rows 1 and 2, respectively, of Table 19-9. (See the earlier section "Comparing groups with conditional distributions" for details.)

To get the numbers in Table 19-9, I started with Table 19-8 and divided the number in each cell by its marginal row total to get a proportion. Each row in Table 19-9 sums to 1 because each row represents its own conditional distribution. (If you're male, you either voted for the incumbent or you didn't — same for females.)

Row 1 of Table 19-9 shows the conditional distribution of voting pattern for males. You see 40% voted for the incumbent and 60% not. Similarly, row 2 of the table shows the conditional distribution of voting pattern for females; again, 40% voted for the incumbent

and 60% did not. Because these distributions are the same, men and women voted the same way; gender and voting pattern are independent.

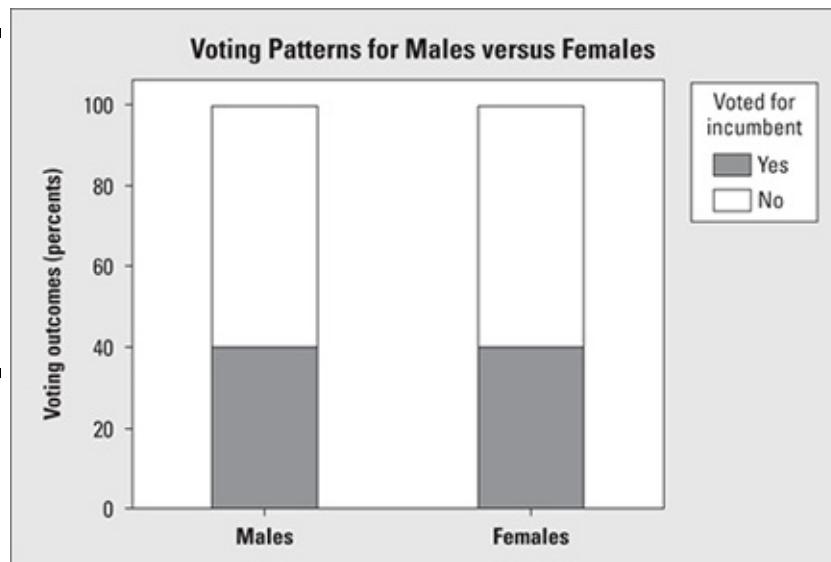
Table 19-9

**Results of Election Survey with
Conditional Distributions**

	<i>Voted for Incumbent President</i>	<i>Didn't Vote for Incumbent President</i>	<i>Total</i>
Males	$44 \div 110 = 0.40$	$66 \div 110 = 0.60$	1.00
Females	$36 \div 90 = 0.40$	$54 \div 90 = 0.60$	1.00

Figure 19-4 shows the conditional distributions of voting pattern for males and females using a graph called a stacked bar chart. Because the bars look exactly alike, you conclude that gender and voting pattern are independent.

Figure 19-4:
Bar graph showing the conditional distributions of voting pattern for males versus females.



To have independence, you don't need the percentages within each bar to be 50-50 (for example, 50% males in favor and 50% males opposed). It's not the percentages within each bar (group) that have to be the same; it's the percentages across the bars (groups) that need to match (for example, 60% of males in favor and 60% of females in favor).



Instead of comparing rows of a two-way table to determine independence, you can compare the columns. In the voting example you'd be comparing the gender breakdowns for the group who voted for the incumbent to the gender breakdowns for the group who didn't vote for the incumbent. The conclusion of independence would be the same as what you found previously, although the percentages you calculate would be different.

Comparing marginal and conditional to check for independence

Another way to check for independence is to see whether the marginal distribution of voting pattern (overall) equals the conditional distribution of voting pattern for each of the gender groups (males and females). If these distributions are equal, then gender doesn't matter. Again, gender and voting pattern are independent.

Looking at the voting pattern example, you find the conditional distribution of voting pattern for the males (first bar in Figure 19-4) is 40% yes and 60% no. To find the marginal (overall) distribution of voting pattern (males and females together), take the marginal column totals in the last row of Table 19-8 (80 yes and 120 no) and divide through by 200 (the grand total). You get $80 \div 200 = 0.40$ or 40% yes, and $120 \div 200 = 0.60$ or 60% no. (See the section "Calculating marginal distributions" earlier in this chapter for more explanation.) The marginal distribution of overall voting pattern matches the conditional distribution of voting pattern for males, so voting pattern is independent of gender.



Here's where a small table with only two rows and two columns cuts you a break. You have to compare only one of the conditionals to the marginal because you have only two groups to compare. If the voting pattern for the males is the same as the overall voting pattern, then the same will be true for the females. To check for independence when you have more than two groups, you use a Chi-square test (discussed in my book *Statistics II For Dummies*, published by Wiley).

Describing a dependent relationship

Two categorical variables are *dependent* if the conditional distributions are different for at least two of the groups being compared. In the election example from the previous section, the groups are males and females, and the variable being compared is whether the person voted for the incumbent president.

Dependence in this case means knowing that the outcome of the first variable does affect the outcome of the second variable. In the election example, if dependence had been found, it would mean that males and females didn't have the same voting pattern for the incumbent (for example, more males voting for the incumbent than females). (Pollsters use this kind of data to help steer their campaign strategies.)



Other ways of saying two variables are dependent are to say they are related, or associated. However, statisticians don't use the term *correlation* to indicate relationships between categorical variables. The word *correlation* in this context applies to the linear relationship between two numerical variables (such as height and weight), as seen in Chapter 18. (This mistake occurs in the media all the time, and it drives us statisticians crazy!)

Here's an example to help you better understand dependence: A recent press release put

out by The Ohio State University Medical Center caught my attention. The headline said that aspirin can prevent polyps in colon-cancer patients. Having had a close relative who succumbed to this disease, I was heartened at the prospect that researchers are making progress in this area and decided to look into it.

The researchers studied 635 colon-cancer patients; they randomly assigned approximately half of them to an aspirin regimen (317 people) and the other half to a placebo (fake pill) regimen (318 people). They followed the patients to see which ones developed subsequent polyps and which did not. The data from the study are summarized in Table 19-10.

Table 19-10 Summary of Aspirin and Polyps Study Results

	<i>Developed Subsequent Polyps</i>	<i>Didn't Develop Subsequent Polyps</i>	<i>Total</i>
Aspirin	54 (17%)	263 (83%)	317 (100%)
Placebo	86 (27%)	232 (73%)	318 (100%)
Total	140	495	635

Comparing the results in the rows of Table 19-10 to check for independence means finding the conditional distribution of outcomes (polyps or not) for the aspirin group and comparing it to the conditional distribution of outcomes for the placebo group. Making these calculations, you find that $54 \div 317 = 17\%$ of patients in the aspirin group developed polyps (the rest, 83%, did not), compared to $86 \div 318 = 27\%$ of the placebo group who developed subsequent polyps (the rest, 73%, did not).

Because the percentage of patients developing polyps is much smaller for the aspirin group compared to the placebo group (17% versus 27%), a dependent relationship appears to exist between aspirin-taking and the development of subsequent polyps among the colon-cancer patients in this study. (But does it carry over to the population? You find out in the section “Projecting from sample to population” later in this chapter.)

Cautiously Interpreting Results

It's easy to get carried away when a relationship between two variables has been found; you see this happen all the time in the media. For example, a study reports that eating eggs doesn't affect your cholesterol as once thought; in the details of the report you see the study was conducted on a total of 20 men who were all in excellent health, on low-fat diets, who exercised several times a week. Ten men in good health ate two eggs a day and their cholesterol didn't change much, compared to ten men who didn't eat two eggs per day. Do these results carry over to the entire population? Can't tell — the subjects in the study don't represent the rest of us. (See Chapter 17 for the scoop on evaluating experiments.)

In this section, you see how to put the results from a two-way table into proper perspective in terms of what you can and can't say and why. This basic understanding gives you the ability to critically evaluate and make decisions about results presented to you (not all of which are correct).

Checking for legitimate cause and effect

Researchers studying two variables often look for links that indicate a cause-and-effect relationship. A *cause-and-effect relationship* between two categorical variables means as you change the value of one variable and all else remains the same, it causes a change in the second variable — for example, if being on an aspirin regimen decreases the chance of developing subsequent polyps in colon-cancer patients.

However, just because two variables are found to be related (dependent) doesn't mean they have a cause-and-effect relationship. For example, observing that people who live near power lines are more likely to visit the hospital in a year's time due to illness doesn't necessarily mean the power lines caused the illnesses.



The most effective way to conclude a cause-and-effect relationship is by conducting a well-designed experiment (where possible). All the details are laid out in Chapter 17, but I touch on the main points here. A well-designed experiment meets the following three criteria:

- ✓ It minimizes *bias* (systematic favoritism of subjects or outcomes).
- ✓ It repeats the experiment on enough subjects so the results are reliable and repeatable by another researcher.
- ✓ It controls for other variables that may affect the outcome that weren't included in the study.

In the earlier section “Describing a dependent relationship,” I discuss a study involving the use of aspirin to prevent polyps in cancer patients. Because of the way the data was collected for this study, you can be confident about the conclusions drawn by the researchers; this study was a well-designed experiment, according to the criteria established in Chapter 17. To avoid problems, the researchers in this study did the following:

- ✓ Randomly chose who took the aspirin and who received a fake pill
- ✓ Had large enough sample sizes to obtain accurate information
- ✓ Controlled for other variables by conducting the experiment on patients in similar situations with similar backgrounds

Because their experiment was well-designed, the researchers concluded that a cause-and-effect relationship was found for the patients in this study. The next test is to see whether they can project these results to the population of all colon-cancer patients. If so, they are truly entitled to the headline “Aspirin Prevents Polyps in Colon-Cancer Patients.” The next section walks you through the test.



Whether two related variables are found to be causally associated depends on how the study was conducted. A well-designed experiment is the most convincing way to establish cause and effect. In cases where an experiment would be unethical (for example, proving that smoking causes lung cancer by forcing people to smoke), a mountain of convincing observational studies (where you collect data on people who smoke and people who don't) would be needed to show that an association between two variables crosses over into a cause-and-effect relationship.

Projecting from sample to population

In the aspirin/polyps experiment discussed in the earlier section “Describing a dependent relationship,” I compare the percentage of patients developing subsequent polyps for the aspirin group versus the non-aspirin group and got the results 17% and 27%, respectively. For this sample, the difference is quite large, so I’m cautiously optimistic that these results would carry over to the population of all cancer patients. But what if the numbers were closer, such as 17% and 20%? Or 17% compared to 19%? How different do the proportions have to be in order to signal a meaningful association between the two variables?



Percentages compared using data from your sample reflect relationships within your sample. However, you know that results change from sample to sample. To project these conclusions to the population of all colon-cancer patients (or any population being studied), the difference in percentages found by the sample has to be *statistically significant*. Statistical significance means even though you know results will vary, even taking that variation into account it's very unlikely the differences were due to chance. That way, the same conclusion about a relationship can be made about the whole population, not just for a particular data set.

I analyzed the data from the aspirin/polyps study using a hypothesis test for the difference of two proportions (found in Chapter 15). The proportions being compared were the proportion of patients taking aspirin who developed subsequent polyps and the proportion of patients not taking aspirin who developed subsequent polyps. Looking at these results, my *p*-value is less than 0.0024. (A *p*-value measures how likely you were to have gotten the results from your sample if the populations really had no difference; see Chapter 14 to get the scoop on *p*-values.)

Because this *p*-value is so small, the difference in proportions between the aspirin and non-aspirin groups is declared to be statistically significant, and I conclude that a relationship exists between taking aspirin and developing fewer subsequent polyps.



You can't make conclusions about relationships between variables in a population based only on the sample results in a two-way table. You must take into account the fact that results change from sample to sample. A hypothesis test gives limits for how different the sample results can be to still say the variables are independent. Beware of conclusions based only on sample data from a two-way table.

Making prudent predictions

A common goal of research (especially medical studies) is to make predictions, recommendations, and decisions after a relationship between two categorical variables is found. However, as a consumer of information, you have to be very careful when interpreting results; some studies are better designed than others.

The colon-cancer study from the previous section shows that patients who took aspirin daily had a lower chance of developing subsequent polyps (17% compared to 27% for the non-aspirin group). Because this was a well-designed experiment and the hypothesis test for generalizing to the population was significant, making predictions and recommendations for the population of colon-cancer patients based on these sample results is appropriate. They've indeed earned the headline of their press release: "Aspirin Prevents Polyps in Colon-Cancer Patients."

Resisting the urge to jump to conclusions



Try not to jump to conclusions when you hear or see a relationship being reported regarding two categorical variables. Take a minute to figure out what's really going on, even when the media wants to sweep you away with a dramatic result.

For example, as I write this, a major news network reports that men are 40% more likely to die from cancer than women. If you're a man, you may think you should panic. But when you examine the details, you find something different. Researchers found that men are much less likely to go to the doctor than women, so by the time cancer is found, it's more advanced and difficult to treat. As a result, men were more likely to die of cancer after its diagnosis. (They aren't necessarily more likely to *get* cancer; that's for a different study.) This study was meant to promote early detection as the best protection and encourage men to keep their annual checkups. The message would have been clearer had the media reported it correctly (but that's not as exciting or dramatic).

Part VI

The Part of Tens

The 5th Wave

By Rich Tennant



In this part . . .

Where would a statistics book be without some statistics of its own? This part contains ten methods for being a statistically savvy sleuth and ten tips for boosting your score on a statistics exam. You can use this quick, concise reference to help critique or design a survey, detect common statistical abuses, and ace your introductory statistics course.

Chapter 20

Ten Tips for the Statistically Savvy Sleuth

In This Chapter

- ▶ Recognizing common statistical mistakes made by researchers and the media
 - ▶ Avoiding mistakes when doing your statistics
-

This book is not only about understanding the statistics that you come across in the media and in your workplace; it's even more about digging deeper to examine whether those statistics are correct, reasonable, and fair. You have to be vigilant — and a bit skeptical — to deal with today's information explosion, because many of the statistics you find are wrong or misleading, either by error or by design. If you don't critique the information you're consuming, in terms of its correctness, completeness, and fairness, who will? In this chapter, I outline ten tips for detecting common statistical mistakes made by researchers and by the media and ways to avoid making them yourself.

Pinpoint Misleading Graphs

Most graphs and charts contain great information that makes a point clearly, concisely, and fairly. However, many graphs give incorrect, mislabeled, and/or misleading information; or they simply lack important information that the reader needs to make critical decisions about what is being presented. Some of these shortcomings occur by mistake; others are incorporated by design in hopes you won't notice. If you're able to pick out problems with a graph before you contemplate any conclusions, you won't be taken in by misleading graphs.

Figure 20-1 shows examples of four important types of data displays: pie charts, bar graphs, time charts, and histograms. In this section I point out some of the ways you can be misled if these types of graphs are not made properly. (For more information on making charts and graphs correctly and identifying misleading ones, see Chapters 6 and 7.)

Pie charts

Pie charts are exactly what they sound like: circular (pie-shaped) charts that are divided into slices that represent the percentage (relative frequency) of individuals that fall into different groups. Groups represent a categorical variable, such as gender, political party, or employment status. Figure 20-1a is a pie chart showing a breakdown of voter opinions

on some issue (call it Issue 1).

Here's how to sink your teeth into a pie chart and test it for quality:

- ✓ Check to be sure the percentages add up to 100 percent, or close to it (any round-off error should be small).
- ✓ Be careful when you see a slice of the pie called "other"; this is the catch-all category. If the slice for "other" is too large (larger than other slices), the pie chart is too vague. On the other extreme, pie charts with many tiny slices give you information overload.
- ✓ Watch for distortions that come with the three-dimensional ("exploded") pie charts, in which the slice closest to you looks larger than it really is because of the angle at which it's presented.
- ✓ Look for a reported total number of individuals who make up the pie chart so you can determine how big the sample was before it was divided up into slices. If the size of the data set (the number of respondents) is too small, the information isn't reliable.

Bar graphs

A bar graph is similar to a pie chart, except that instead of being in the shape of a circle that's divided up into slices, a bar graph represents each group as a bar, and the height of the bar represents the number (frequency) or percentage (relative frequency) of individuals in that group. Figure 20-1b is a relative frequency–style bar graph showing voter opinions on some issue (call it Issue 1); its results correspond with the pie chart shown in Figure 20-1a.

When examining a bar graph:

- ✓ Check for the sample size. If the bars represent frequencies, you find the sample size by summing them up; if the bars represent relative frequencies, you need the sample size to know how much data went into making the graph.
- ✓ Consider the units being represented by the height of the bars and what the results mean in terms of those units. For example, are they showing the total number of crimes, or the crime rate (also known as total number of crimes per capita)?
- ✓ Evaluate the starting point of the axis where the counts (or percents) are shown, and watch for the extremes: If the heights of the bars fluctuate from 200 to 300 but the counts axis starts at 0, the heights of the bars won't look much different. However, if the starting point on the counts axis is 200, you are basically chopping off the bottoms of all the bars, and what differences remain (ranging

from 0 to 100) will look more dramatic than they should.

- ✓ Check out the range of values on the axis where the counts (or percents) are shown. If the heights of the bars range from 6 to 108 but the axis shows 0 to 500, the graph will have a great deal of white space and differences in the bars become hard to distinguish. However, if the axis goes from 5 to 110 with almost no breathing room, the bars will be stretched to the limit, making differences between groups look larger than they should.

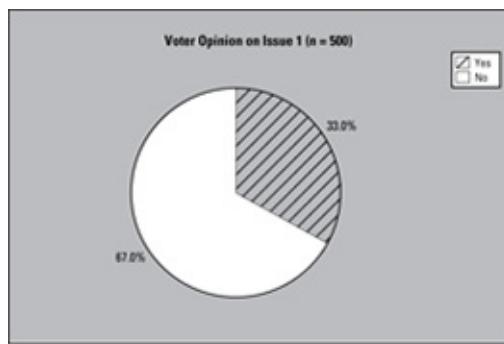
Time charts

A time chart shows how a numerical variable changes over time (for example, stock prices, car sales, or average temperature). Figure 20-1c is an example of a time chart showing the percentage of yes voters from 2002 to 2010, in 2-year increments.

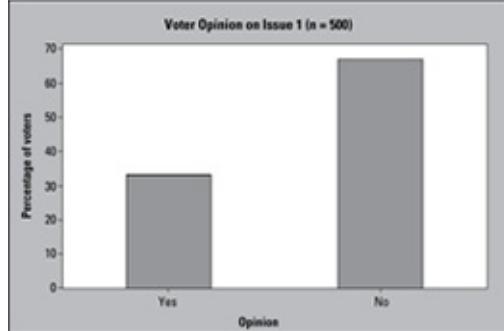
Here are some issues to watch for with time charts:

- ✓ Watch the scale on the vertical (quantity) axis as well as the horizontal (timeline) axis; results can be made to look more or less dramatic than they actually are by simply changing the scale.
- ✓ Take into account the units being portrayed by the chart and be sure they are equitable for comparison over time; for example, are dollar amounts being adjusted for inflation?
- ✓ Beware of people trying to explain why a trend is occurring without additional statistics to back themselves up. A time chart generally shows *what* is happening. *Why* it's happening is another story!
- ✓ Watch for situations in which the time axis isn't marked with equally spaced jumps. This often happens when data are missing. For example, the time axis may have equal spacing between 2001, 2002, 2005, 2006, 2008 when it should actually show empty spaces for the years in which no data are available.

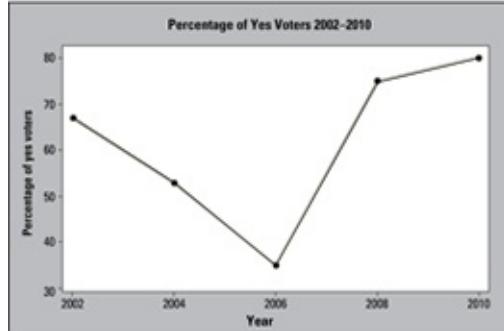
Figure 20-1:
Four types of graphs: a) pie chart; b) bar graph; c) time chart; and d) histogram.



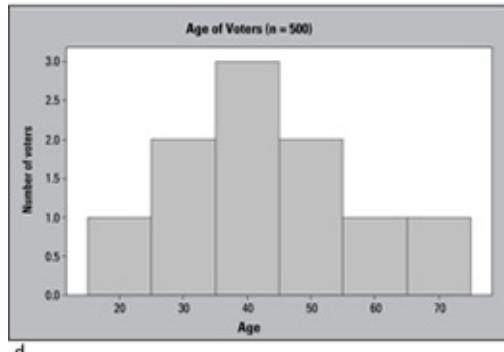
a



b



c



d

Histograms

A *histogram* is a graph that breaks the sample into groups according to a numerical variable (such as age, height, weight, or income) and shows either the number of individuals (frequency) or the percentage of individuals (relative frequency) that fall into each group. Figure 20-1d is a frequency style histogram showing the ages of voters in a certain election.

Some items to watch for regarding histograms include the following:

- ✓ Watch the scale used for the vertical (frequency/relative frequency) axis, looking

especially for results that are exaggerated or played down through the use of inappropriate scales.

- ✓ Check out the units on the vertical axis to see whether they report frequencies or relative frequencies; if they're relative frequencies, you need the sample size to determine how much data you're looking at.
- ✓ Look at the scale used for the groupings of the numerical variable on the horizontal axis. If the groups are based on small intervals (for example, 0–2, 2–4, and so on), the heights of the bars may look choppy and overly volatile. If the groups are based on large intervals (for example, 0–100, 100–200, and so on), the data may give a smoother appearance than is realistic.

Uncover Biased Data

Bias in statistics is the result of a systematic error that either overestimates or underestimates the true value. For example, if I use a ruler to measure plants and that ruler is 1/2-inch short, all of my results are biased; they're systematically lower than their true values.

Here are some of the most common sources of biased data:

- ✓ Measurement instruments may be systematically off. For example, a police officer's radar gun may say you were going 76 miles per hour when you know you were only going 72 miles per hour. Or a badly adjusted scale may always add 5 pounds to your weight.
- ✓ The way the study is designed can create bias. For example, a survey question that asks, "Have you *ever* disagreed with the government?" will overestimate the percentage of people who are generally unhappy with the government. (See Chapter 16 for ways to minimize bias in surveys.)
- ✓ The sample of individuals may not represent the population of interest — for example, examining student study habits by only going to the campus library. (See more in the section, "Identify Non-Random Samples" later in this chapter.)
- ✓ Researchers aren't always objective. Suppose in a drug study one group of patients is given a sugar pill and the other group is given the real drug. If the researchers know who received the real drug, they may inadvertently pay more attention to those patients to see if it's working; they may even project results onto the patients (such as saying, "I bet you're feeling better, aren't you?"). This creates a bias in favor of the drug. (See Chapter 17 for more information on setting up good experiments.)



To spot biased data, examine how the data were collected. Ask questions about the selection of the participants, how the study was conducted, what questions were used, what treatments (medications, procedures, therapy, and so on) were given (if any) and who knew about them, what measurement instruments were used and how they were calibrated, and so on. Look for systematic errors or favoritism, and if you see too much of it, ignore the results.

Search for a Margin of Error

The word *error* has a somewhat negative connotation, as if an error is something that is always avoidable. In statistics, that's not always the case. For example, a certain amount of what statisticians call *sampling error* will always occur whenever someone tries to estimate a population value using anything other than the entire population. Just the act of selecting a sample from the population means you leave out certain individuals, and that means you're not going to get the precise, exact population value. No worries, though. Remember that statistics means never having to say you're certain — you have to only get close. And if the sample is large enough, the sampling error will be small (assuming it's good data of course).

To evaluate a statistical result, you need a measure of its accuracy — typically through the margin of error. The margin of error tells you how much the researcher expects her results to vary from sample to sample. (For more information on margin of error, see Chapter 12.) When a researcher or the media fail to report the margin of error, you're left to wonder about the accuracy of the results, or worse, you just assume that everything is fine, when in many cases, it's not.



When looking at statistical results in which a number is being estimated (for example, the percentage of all Americans who think the president is doing a good job), always check for the margin of error. If it's not included, ask for it! (Or if given enough other pertinent information, you can calculate the margin of error yourself using the formulas in Chapter 13.)

Identify Non-Random Samples

If you're trying to study a population but you can only study a sample of individuals from it, how can you ensure that your sample represents the population? The most important criteria is to select your sample in a random fashion; that is, to take a *random sample*. You know a sample is random if it had the same chance of being selected as every other possible sample of the same size did. (It's like pulling names from a hat.)

Many surveys aren't based on random samples, however. For example, TV polls asking viewers to "call us with your opinion" don't represent random samples. In fact they don't represent samples at all; when you take a sample, you select individuals from the population; for call-in polls, the individuals select themselves.

Experiments (particularly medical studies) typically can't involve a random sample of individuals, for ethical reasons. You can't call someone and say, "You were chosen at random to participate in a sleep study. You'll need to come down to our lab tomorrow and stay there for two nights." Such types of experiments are conducted using subjects that volunteer to participate — they're not randomly selected first.



But even though you can't randomly select the subjects (participants) for your experiment, you can still get valid results if you incorporate the randomness in a different way — by randomly assigning the subjects to the treatment group and the control group. If the groups were assigned at random, they have a good chance of being very similar, except for what treatment they received. That way, if you do find a large enough difference in the outcomes of the groups, you can attribute those differences to the treatment, rather than to other factors.



Before making any decisions about statistical results from a survey, look to see how the sample of individuals was selected. If the sample wasn't selected randomly, take the results with a grain of salt (see Chapter 16). If you're looking at the results of an experiment, find out whether the subjects were randomly assigned to the treatment and control groups; if not, ignore the results (see Chapter 17).

Sniff Out Missing Sample Sizes

Both the quality and quantity of information is important in assessing how accurate a statistic will be. The more good data that goes into a statistic, the more accurate that statistic will be. The quality issue is tackled in the section "Uncover Biased Data" earlier in this chapter. When the quality has been established, you need to assess the accuracy of the information, and for that you need to look at how much information was collected (that is, you have to know the sample size).

Small sample sizes make results less accurate (unless your population was small to begin with). Many headlines aren't exactly what they appear to be when the details reveal a study that was based on a small sample. Perhaps even worse, many studies don't even report the sample size at all, which should lead you to be skeptical of the results. (For example, an old chewing gum ad said, "Four out of five dentists surveyed recommend [this gum] for their patients who chew gum." What if they really did ask only five dentists?)



Don't think about this too much, but according to statisticians (who are picky about precision), 4 out of 5 is much different than 4,000 out of 5,000, even though both fractions equal 80 percent. The latter represents a much more precise (repeatable) result because it's based on a much higher sample size. (Assuming it's good data, of course.) If you ever wondered how math and statistics are different, here's your answer! (Chapter 12 has more on precision.)

However, more data isn't always better data — it depends on how well the data were collected (see Chapter 16). Suppose you want to gather the opinions of city residents on a city council proposal. A small random sample with well-collected data (such as a mail survey of a small number of homes chosen at random from a city map) is much better than a large non-random sample with poorly collected data (for example, posting a Web survey on the city manager's Web site and asking for people to respond).



Always look for the sample size before making decisions about statistical information. The smaller the sample size, the less precise the information. If the sample size is missing from the article, get a copy of the full report of the study, contact the researcher, or contact the journalist who wrote the article.

Detect Misinterpreted Correlations

Everyone wants to look for connections between variables; for example, what age group is more likely to vote Democrat? If I take even more vitamin C, am I even less likely to get a cold? How does staring at the computer all day affect my eyesight? When you think of connections or associations between variables, you probably think of correlation. Yes, correlation is one of the most commonly used statistics — but it's also one of the most misunderstood and misused, especially throughout the media.

Some important points about correlation are as follows (see Chapter 18 for all the additional information):

- ✓ **The statistical definition of *correlation* (denoted by r) is the measure of strength and direction of the linear relationship between two numerical variables.** A correlation tells you whether the variables increase together or go in opposite directions and the extent to which the pattern is consistent across the data set.
- ✓ **The statistical term *correlation* is only used in the context of two numerical variables (such as height and weight).** It does not apply to two categorical variables (such as political party and gender).

For example, voting pattern and gender may be related, but using the word

correlated to describe their relationship isn't "sc" (statistically correct, get it?). You can say two categorical variables are *associated*.

- ✓ **If a strong correlation and scatterplot exist between two numerical variables, you should be able to draw a straight line through the points, and the points should lie close to the line.** If a line doesn't fit the data well, the variables likely won't have a strong correlation (r), and vice versa. (See Chapter 18 for information on line-fitting, also known as *linear regression*.)



A weak correlation implies that a linear relationship doesn't exist between the two variables, but this doesn't necessarily mean the variables aren't related at all. They may have some other type of relationship besides a linear relationship. For example, bacteria multiply at an exponential rate over time (their numbers explode, doubling faster and faster).

- ✓ **Correlation doesn't automatically mean cause and effect.** For example, suppose Susan reports based on her observations that people who drink diet soda have more acne than people who don't. If you're a diet soda drinker, don't break out just yet! This correlation may be a freak coincidence that only happened to the people she observed. At most, it means more research needs to be done (beyond observation) in order to draw any connections between diet soda and acne. (Susan can read Chapter 17 to find out how to design a good experiment.)

Reveal Confounding Variables

A *confounding variable* is a variable that isn't included in a study but whose influence can affect the results and create confusing (confounding) conclusions. For example, suppose a researcher reports that eating seaweed helps you live longer, but when you examine the study, you find out that it was based on a sample of people who regularly eat seaweed in their diets and are over the age of 100. When you read the interviews of these people, you discover some of their other secrets to long life (besides eating seaweed): They slept an average of 8 hours a day, drank a lot of water, and exercised every day. So did the seaweed cause them to live longer? You can't tell, because several confounding variables (exercise, water consumption, and sleeping patterns) may also have contributed.



The best way to control for confounding variables is to conduct a well-designed experiment (see Chapter 17), which involves setting up two groups that are alike in as many ways as possible, except that one group receives a specified treatment and the other group receives a control (a fake treatment, no treatment, or a standard, non-experimental treatment). You then compare the results from the two groups, attributing any significant differences to the treatment (and to nothing else, in an ideal world).



This seaweed study wasn't a designed experiment; it was an observational study. In observational studies, no control for any variables exists; people are merely observed, and information is recorded. Observational studies are great for surveys and polls, but not for showing cause-and-effect relationships, because they don't control for confounding variables. A well-designed experiment provides much stronger evidence.

If doing an experiment is unethical (for example, showing smoking causes lung cancer by forcing half of the subjects in the experiment to smoke ten packs a day for 20 years while the other half of the subjects smoke nothing), then you must rely on mounting evidence from many observational studies over many different situations, all leading to the same result. (See Chapter 17 for all the details on designing experiments.)

Inspect the Numbers

Just because a statistic appears in the media doesn't mean it's correct. In fact, errors appear all the time (by mistake or by design), so stay on the lookout for them. Here are some tips for spotting botched numbers:

- ✓ **Make sure everything adds up to what it's reported to.** With pie charts, be sure all the percentages add up to 100 percent (subject to a small amount of rounding error).
- ✓ **Double-check even the most basic of calculations.** For example, a pie chart shows that about 83.33 percent of Americans are in favor of an issue, but the accompanying article reports "7 out of every 8" Americans are in favor of the issue. Are these statements saying the same thing? No; 7 divided by 8 is 87.5 percent — if you want 83.33 percent, it's 5 out of 6.
- ✓ **Look for the response rate of a survey; don't just be happy with the number of participants.** (The response rate is the number of people who responded divided by the total number of people surveyed times 100 percent.) If the response rate is much lower than 50 percent, the results may be biased, because you don't know what the non-respondents would have said. (See Chapter 16 for the full scoop on surveys and their response rates.)
- ✓ **Question the type of statistic used, to determine whether it's appropriate.** For example, suppose the number of crimes went up, but so did the population size. Instead of reporting the number of crimes, the media need to report the crime rate (number of crimes per capita).



Statistics are based on formulas that take the numbers you give them and crunch

out what you ask them to crunch out. The formulas don't know whether the final answers are correct or not. The people behind the formulas should know better, of course. Those who don't know better will make mistakes; those who do know better might fudge the numbers anyway and hope you don't catch on. You, as a consumer of information (also known as a certified skeptic), must be the one to take action. The best policy is to ask questions.

Report Selective Reporting

You cannot credit studies in which a researcher reports his one statistically significant result but fails to mention the reports of his other 25 analyses, none of which came up significant. If you had known about all the other analyses, you may have wondered whether this one statistically significant result is truly meaningful, or simply due to chance (like the idea that a monkey typing randomly on the typewriter would eventually write Shakespeare). It's a legitimate question.

The misleading practice of analyzing data until you find something is what statisticians call *data snooping* or *data fishing*. Here's an example: Suppose Researcher Bob wants to figure out what causes first graders to argue with each other so much in school (he must not be a parent or he wouldn't even try to touch this one!). He sets up a study in which he observes a classroom of first graders every day for a month and records their every move. He gets back to his office, enters all his data, hits a button that asks the computer to perform every analysis known to man, and sits back in his chair eagerly awaiting the results. After all, with all this data he's bound to find *something*.

After poring through his results for several days, he hits pay dirt. He runs out of his office and tells his boss he's got to put out a press release saying a ground-breaking study finds that first graders argue most when 1) the day of the week ends in the letter y or 2) when the goldfish in their classroom aquarium swims through the hole in its sunken pirate ship. Great job, Researcher Bob! I've got a feeling that a month of watching a group of first graders took the edge off his data analysis skills.



The bottom line is that if you collect enough data and analyze it long enough, you're bound to find something, but that something may be totally meaningless or just a fluke that's not repeatable by other researchers.

How do you protect yourself against misleading results due to data fishing? Find out more details about the study, starting with how many tests were done in total, and how many of those tests were found to be non-significant. In other words, get the whole story if you can, so that you can put the significant results into perspective.



To avoid being reeled in by someone's data fishing, don't just go with the first result that you hear, especially if it makes big news and/or seems a little suspicious. Contact the researchers and ask for more information about their data, or wait to see whether other researchers can verify and replicate their results.

Expose the Anecdote

Ah, the anecdote — one of the strongest influences on public opinion and behavior ever created. And one of the least valid. An *anecdote* is a story or result based on a single person's experience or situation. For example:

- ✓ The waitress who won the lottery — twice.
- ✓ The cat that learned how to ride a bicycle.
- ✓ The woman who lost a hundred pounds in two days on the new miracle potato diet.
- ✓ The celebrity who claims to have used an over-the-counter hair color for which she is a spokesperson (yeah, right).

Anecdotes make great news; the more sensational the better. But sensational stories are outliers from the norm of life. They don't happen to most people.

You may think you're out of reach of the influence of anecdotes. But what about those times when you let one person's experience influence you? Your neighbor loves his Internet service provider, so you try it, too. Your friend had a bad experience with a certain brand of car, so you don't bother to test-drive it. Your dad knows somebody who died in a car crash because she was trapped in the car by her seat belt, so he decides never to wear his.

While some decisions are okay to make based on anecdotes, some of the more important decisions you make should be based on real statistics and real data that come from well-designed studies and careful research.



An anecdote is really a data set with a sample size of only one. You have no information to compare it to, no statistics to analyze, no possible explanations or information to go on — just a single story. Don't let anecdotes have much influence over you. Instead, rely on scientific studies and statistical information based on large random samples of individuals who represent their target populations (not just a single situation). When someone tries to persuade you by telling you an anecdote just say, "Show me the data!"

Chapter 21

Ten Surefire Exam Score Boosters

In This Chapter

- ▶ Getting into the zone
 - ▶ Developing savvy strategies
 - ▶ Preventing silly mistakes
-

I've taught more than 40,000 students in my teaching career (don't try to guess how old I am, it's not polite!), and each student has taken at least three exams for me. That makes over 120,000 exams I've graded or had a hand in grading, and believe me, I've seen it all. I've seen excellent answers, disastrous answers, and everything in between. I've gotten notes from students in the margins asking me to go easy on them because their dog ran away and they didn't have time to study. I've seen some answers that even I couldn't figure out. I've laughed, I've cried, and I've beamed with pride at what my students have come up with in exam situations.

In this chapter, I've put together a list of ten strategies most often used by students who do well on exams. These students are not necessarily smarter than everyone else (although you do have to know your material, of course), but they are much better prepared. As a result, they are able to handle new problems and situations without getting thrown off; they make fewer little mistakes that chip away at an exam score; and they are less likely to have that deer-in-the headlights look, not being able to start a problem. They are more likely to get the right answer (or at least get partial credit) because they are good at labeling information and organizing their work. No doubt about it — preparation is the key to success on a stat exam.

You too can be a successful statistics student — or *more* successful, if you're already doing well — by following the simple strategies outlined in this chapter. Remember, every point counts, and they all add up, so let's start boosting your exam score right away!

Know What You Don't Know, and then Do Something about It

Figuring out what you know and what you don't know can be hard when you are taking a statistics class. You read the book and can understand all the examples in your notes,

but you can't do your homework problems. You can answer all your roommate's statistics questions, but you can't answer your own. You walk out of an exam thinking you did well, but when you see your grade, you are shocked.

What's happening here? The bottom line is, you have to be aware of what you know and what you don't know if you want to be successful. This is a very tough skill to develop, but it's well worth it. Students often find out what they don't know the hard way — by losing points on exam questions. Mistakes are okay, we all make them — what matters is *when* you make them. If you make a mistake before the exam while you still have time to figure out what you're doing wrong, it doesn't cost you anything. If you make that same mistake on an exam, it'll cost you points.



Here's a strategy for figuring out what you know and what you don't know. Go through your lecture notes and place stars by any items from the notes that you don't understand. You can also "test" yourself, as I describe later in "Yeah-yeah trap #2," and make a list of problems that stumped you. Take your notes and list to your professor and ask him to go through the problem areas with you. Your questions will be specific enough that your professor can zoom in when he's talking with you, give you specific information and examples, and then check to make sure you understand each idea before moving on to the next item. Meeting with your professor won't take long; sometimes getting one question answered has a ripple effect and clears up other questions farther down on your list.



Leave no stone unturned when it comes to making sure you understand all the concepts, examples, formulas, notation, and homework problems before you walk into the exam. I always tell my students that 30 minutes with me has a potential of raising your grade by 10%, because I'm awfully good at explaining things and answering questions — and I'm probably better at it than any roommate, brother-in-law, or friend who took the class four years ago with another professor. A quick office visit with your professor is well worth your time — especially if you bring a detailed list of questions with you. If for some reason your professor is not available, see if you have access to a tutor for help.

All-purpose pointers for succeeding in class

Here's some general advice my students have found helpful:

- ✓ I know you've heard this before, but you really are at an advantage if you go to class every day so you have a full set of notes to review. It also ensures you didn't miss any of the little things that add up to big points on an exam.
- ✓ Don't just write down what the professor wrote down — that's for amateurs. The professionals also write down anything else he made a big deal about but didn't write down. That's what separates

the As from the Bs.

✓ Do little things to stay organized while you go through the course; you won't get overwhelmed later when it's crunch time. The day I invested 5 dollars and bought a good mechanical pencil, a good eraser, a cheap three-hole punch for my handouts, and a tiny stapler was one of the best days of my student life. Okay, it'll probably cost you 10 dollars for these items today, but trust me, it'll be worth it!

✓ Get to know your professor and let her get to know you. Introducing yourself on the first day makes a big impression; getting face time (as well as some good help) by asking a question after class (if you have one) or stopping in during office hours never hurts. Don't worry about whether your questions are silly — it's not what level you're at now that counts; it's your desire to get to the next level and do well in the class that's important. That's what your professor wants to see.

Avoid “Yeah-Yeah” Traps

What's a “yeah-yeah” trap? It's a term I use when you get caught saying “Yeah-yeah, I got this; I know this, no problem,” but then comes the exam and whoa — you didn't have it, you didn't know it, and Houston, you actually had a problem. Yeah-yeah traps are bad because they lull you into thinking you know everything, you don't have any questions, and you'll get 100% on the exam, when the truth is you still need to resolve some issues.

Although many different yeah-yeah traps exist, I point out the two most common ones in this section and help you avoid them. I call them (cleverly) *yeah-yeah trap #1* and *yeah-yeah trap #2*. Both of these traps are subtle, and they can sneak up on even the most conscientious students, so if you recognize yourself in this section, don't feel bad. Just think how many points you'll be saving yourself when you get out of “yeah-yeah” mode and into “wait a minute — here's something I need to get straightened out!” mode.

Yeah-yeah trap #1

Yeah-yeah trap #1 happens when you study by looking through your lecture notes over and over again, saying “yeah, I get that,” “I understand that,” and “okay, I can do that,” but you don't actually try the problems from scratch totally on your own. If you understand a problem that's already been done by someone else, it only means you understand what that person did when *they* worked the problem. It doesn't say anything about whether you could have done it on your own in an exam situation when the pressure is on and you're staring at a blank space where your answer is supposed to be. Big difference!

I fall into yeah-yeah trap #1 too. I read through my DVR (digital video recording) manual from beginning to end, and it all made total sense to me. But a week later when I went to record a movie, I had no clue how to do it. Why not? I understood the information as I

was reading along, but I didn't try to apply it for myself, and when the time came I couldn't remember how to do it.

Students always tell me, "If someone sets up the problem for me, I can always figure it out." The problem is, almost anyone can solve a problem that's already been set up. In fact, the whole point is being able to set it up, and no one is going to do that for you on an exam.



Avoid yeah-yeah trap #1 by going through your notes, pulling out a set of examples that your professor used, and writing each one on a separate piece of paper (just the problem, not the solution). Then mix up the papers and make an "exam" out of them. For each problem, try to start it by writing down just the very first step. Don't worry about finishing the problems; just concentrate on starting them. After you've done this step for all the problems, go back into your lecture notes and see if you started them right. (On the back of each problem, write down where it came from in your notes so you can check your answers faster.)

Yeah-yeah trap #2

Yeah-yeah trap #2 is even more subtle than yeah-yeah trap #1. A student comes into my office after the exam and says, "Well I worked every problem in the notes, I redid all the homework problems, I worked all the old exams you posted, and I did great on all of them; I hardly got a single problem wrong. But when I took the exam, I bombed it."

What happened? Nine out of ten times, students in yeah-yeah trap #2 did indeed work all those problems, and spent hours upon hours doing so. But whenever they got stuck and couldn't finish a problem, they peeked at the solutions (which they kept sitting right next to them), saw where they went wrong, said "yeah-yeah, that was a silly mistake — I knew that!" and continued on to finish the problem. In the end they thought they got the problems correct all by themselves, but on an exam they lost some (if not all) of the points, depending on where they originally got stuck.

So how do you avoid yeah-yeah trap #2? By making a test run under "real" exam conditions where the pressure is on. Here's how:

- 1. Study as much as you need to, in whatever manner you need to, until you are ready to test your knowledge.**
- 2. Sit down with a practice exam, or if one isn't available, make your own by choosing some problems from homework, your notes, or the book and shuffling them up.**

Just like at a real exam, you also need a pencil, a calculator, and any other materials you are allowed to bring to your exam — and nothing else! Putting your book and notes away may make you feel anxious, frustrated, or exposed when you do a test run

of an exam, but you really need to find out what you can do on your own before you do the real thing.



Some teachers allow you to bring a *review sheet* (also sometimes called a *memory sheet* or — cringe — a *cheat sheet*), a sheet of paper on which you can write any helpful information you want, subject to limitations that your professor may give. If your teacher allows review sheets at tests, use one for your practice test, too.

3. Turn on the oven timer for however long your exam is scheduled to last, and then get started.

4. Work as many problems as you can to the best of your ability, and when you are finished (or time runs out), put your pencil down.

5. When your “exam” is over, get into the lotus position and breathe in, hold it, and breathe out three times. Then look at the solutions and grade your paper the way your professor would.

If you couldn't start a problem, even if you just forgot one little thing and you immediately recognized it when you saw the solutions — you can't say “Yeah-yeah, I knew that; I wouldn't make that mistake on a real exam”; you have to say “No, I couldn't start it on my own. I would have gotten 0 points for that problem. I need to figure this out.”



You don't get a second chance on a real exam, so when you're studying, don't be afraid to admit when you can't do a problem correctly on your own; just be glad you caught it, and figure out how to fix the problem so you'll get it right next time. Go back over it in your notes, read about it in the book, ask your professor, try more problems of the same type, or ask your study buddy to quiz you on it. Also, try to see a pattern in the type of problems that you were missing points on or getting wrong altogether. Figure out why you missed what you missed. Did you read the questions too fast, which caused you to answer them incorrectly? Was it a vocabulary or a notation issue? How did your studying align with what was on the test? And so on.



Being critical of yourself is hard, and finding out you didn't know something you thought you knew is a little scary. But if you put yourself out there and find your mistakes before they cost you points, you'll zoom in on your weaknesses, turn them into strengths, boost your knowledge, and get a higher exam score.

Make Friends with Formulas

Many students are not comfortable with formulas (unless you are a math nerd, in which

case formulas make you shout for joy). That unease is understandable — I used to be intimidated by them too (formulas, that is — not math nerds). The trouble is, you really can't survive too long without eventually using a formula in a statistics class, so becoming comfortable with them right from the start is important. A formula tells you much more than how to calculate something. It shows the thinking process behind the calculations. For example, the big picture regarding standard deviation can be seen by analyzing its formula:

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

Subtracting the mean, \bar{x} , from a value in the data set, x_i , measures how far above or below the mean that number is. Because you don't want the positive and negative differences to cancel each other out, you square them all to make them positive (but remember that this gives you square units). Then you add them up and divide by $n - 1$, which is near to finding an average, and take the square root to get back into original units. In a general sense, you are finding something like the average distance from the mean.

Stepping back even further, you can tell from the formula that the standard deviation can't be negative, because everything is squared. You also know the smallest it can be is zero, which occurs when all the data are the same (that is, all are equal to the mean). And you see how data that is far from the mean will contribute a larger number to the standard deviation than data that is close to the mean.

And here's another perk. Because you understand the formula for standard deviation now, you know what it's really measuring: the spread of the data around the mean. So when you get an exam question saying "Measure the spread around the mean," you'll know what to do. Bam!



In order to feel comfortable about formulas, follow these tips:

- ✓ **Get into the right mind-set.** Think of formulas as mathematical shorthand and nothing more. All you have to do is be able to decipher them. Oftentimes you're allowed to bring a review sheet to your exam, or you'll be given a formula sheet with your exam, so you may not have to make things harder by memorizing them.
- ✓ **Understand every part of every formula.** In order for any formula to be useful, you have to understand all its components. For example, before you can use the formula for standard deviation, you need to know what x_i and \bar{x} mean and what $\sum_{i=1}^n$ stands for. Otherwise it's totally useless.
- ✓ **Practice using formulas from day one.** Use them to verify the calculations done in lecture or in your book. If you get a different answer from what's shown, figure out what you are doing wrong. Making mistakes here is okay — you caught the

problem early, and that's all that counts.

- ✓ **Whenever you use a formula to do a problem, write it down first and then plug in the numbers in the second step.** The more often you write down a formula, the more comfortable you will be using it on an exam. And if (heaven forbid!) you copy the formula down wrong, your instructor will be able to follow your error, which may mean some partial credit for you!



Chances are, if you've learned some formulas in your class, you're going to need to use them on your exam. Don't expect to be able to use formulas with confidence on an exam if you haven't practiced with them and written them down many, many times beforehand. Practice when the problems are easy so when they get harder you won't have to worry as much.

Make an “If-Then-How” Chart

Quarterbacks always talk about trying to get the game to “slow down” for them so they feel like they have more time to think and react. You want the same thing when you take a statistics exam. (See, you and your NFL hero really do have something in common!) The game starts slowing down for a quarterback when he begins to see patterns in the way the defense lines up against him, rather than feeling like every play brings a completely different look. Similarly for you, the exam starts to “slow down” when the problems start falling into categories as you read them, rather than each one appearing to be totally different from anything you've ever seen before.

To make this happen, many of my students find help in making what I call an *if-then-how chart*. An *if-then-how chart* maps out the types of problems you are likely to run into, strategies to solve them, and examples for quick reference. The basic idea of the *if-then-how chart* is to say “*If* the problem asks for X, *then* I solve it by doing Y, and here's how.” An *if-then-how chart* contains three columns:

- ✓ **If:** In the *if* column, write down a succinct description of what you are asked to find or do. For example, if the problem asks you to test a claim about the population mean (see Chapter 14 for more about claims), write “Test a claim — population mean.” If you are asked to give your best estimate of the population mean (Chapter 13 has the scoop on estimates), write “Estimate population mean.”



Problems are worded in different ways, because that's how the real world works. Pay attention to different wordings that in essence boil down to the same problem, and add them to the appropriate place in the *if* column where the actual problem is already listed. For example, one problem may ask you to estimate the population mean; another problem may say, “Give a range of likely values for the

population mean.” These questions ask for the same thing, so include both in your *if* column.

- ✓ **Then:** In your *then* column you write the exact statistical procedure, formula, or technique you need to solve that type of problem using the statistical lingo. For example, when your *if* column says “Test a claim — population mean,” your *then* column should say “Hypothesis test for μ .” When your *if* statement reads “Estimate population mean” your *then* column should read “Confidence interval for μ .”



To match strategies to situations, look carefully at how the examples in your lecture notes and your book were done and use them as your guide.

- ✓ **How:** In the *how* column, write an example, a formula, and/or a quick note to yourself that will spark your mind and send you off running in the right direction. Write whatever you need to feel comfortable (no one’s going to see it but you, so make it your way!). For example, suppose your *if* column says “Estimate the population mean,” and your *then* column says “Confidence interval — population mean.” In the *how* column, you can write the formula.

Although I just took a lot of time and talking to walk you through it, making an if-then-how chart is much easier done than said. Below is an example of an entry in an if-then-how chart for the confidence interval problem I just laid out.

If

Estimate the population mean (also known as range of likely values)

Then

CI for μ

How

$\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$

Using these three columns, fill in your if-then-how chart with each different type of problem you’ve covered in class. Don’t write down every little example; look for patterns in the problems and boil down the number of scenarios to a doable list.



If-then-how charts should be customized to your needs, so the only way it’s going to work is if you make it yourself. No two people think alike; what works for your friend may not work for you. However, it might be helpful to compare your chart with a friend’s once you are both finished, to see if you’ve left anything out.



If you’re allowed to bring a review sheet to exams, I suggest putting your if-then-how chart on one side. On the other side, write down those little nuggets of information your professor gave you in lecture but didn’t write down. If you aren’t allowed to have a review sheet during the exam, call me crazy, but I’ll argue that you should still make one to study from. Making one really helps you sort out all the

ideas so when you take the exam you'll be much more clear about what to look for and how to set up and solve problems. Lots of students come out of an exam saying they didn't even use their review sheet, and that's when you know you've done a good job putting one together: When it went on the sheet, it went into your mind!

Figure Out What the Question Is Asking

Students often tell me that they don't understand what a problem is asking for. That's the million dollar question, isn't it? And it's not a trivial matter. Oftentimes the actual question is embedded somewhere in the language of the problem; it isn't usually as clear as: "Find the mean of this data set."



For example, a question may ask you to "interpret" a statistical result. What does "interpret" really mean? To most professors the word "interpret" means to explain in words that a nonstatistician would understand.

Suppose you are given some computer output analyzing number of crimes and number of police officers, and you are asked to interpret the correlation between them. First you pick off the number from the output that represents the correlation (say it's -0.85); then you talk about its important features in language that is easy for others to understand. The answer I would like to see on an exam goes something like this: "The correlation between number of police officers and number of crimes is -0.85 ; they have a strong negative linear relationship. As the number of police officers increases, number of crimes decreases."



If you know what the problem is asking for, you have a better chance of actually solving it. You'll gain confidence when you know what you are supposed to do. On the flip side, if you don't know what the problem is asking, even starting it will be very hard. Your anxiety will go up, which can affect your ability to work other problems as well. So how do you boil down a problem to figure out exactly what it's asking for? Here are some tips to follow:

- ✓ **Check the very last sentence of the problem — that's usually where the question is located.** Rather than reading the entire problem a second (and third and fourth) time and getting yourself all worked up, just read it once and then focus on the end of the problem.
- ✓ **Practice boiling down questions ahead of time.** Look at all the examples from your lecture notes, your homework problems, and problems in your textbook and try to figure out what each problem is asking for. Eventually you'll start to see patterns in the way problems are worded, and you'll get better at figuring out

what they are really asking for.

- ✓ **Ask your professor what clues you should look for, and bring example problems with you.** She will be impressed because you are trying to figure out the big picture, and oh, how professors love those “big picture” questions! And after she helps you, you can add those to your if-then-how chart (see “Make an ‘If-Then-How’ Chart”).
- ✓ **Translate the wording of the problem into a statistical statement.** This involves labeling not only what you are given (as discussed in the next section), but also what you want to find.

For example, Professor Barb wants to give 20 percent of her students an A on her statistics exam; your job is to find the cutoff exam score for an A, and this translates to “find the score representing the 80th percentile.”

Label What You’re Given



Many students try to work problems by pushing around numbers that are given in the problem. This approach may work with easy problems, but everyone hits the wall at some point and needs more support to solve harder problems. You’ll benefit from getting into the habit of labeling everything properly — labeling is the critical connection between the *if* column and the *then* column in your if-then-how chart (described earlier in this chapter). You may read a problem and know what you need to do, but without understanding how to use what you’re given in the problem, you won’t be able to solve it correctly. To really understand the numbers the problem gives you, take each one and write down what it stands for.

Suppose you’re given the following problem to solve: “You want to use the size of a house in a certain city (in square feet) to predict its price (in thousands). You collect data on 100 randomly selected homes that have recently been sold. You find the mean price is \$219,100 with standard deviation of \$60,100, and you know the mean size is 1,993 square feet, with standard deviation of 349 square feet. You find the correlation between size and price for these homes is +0.90. Find the best-fitting regression line that you can use to predict house price using size.”

Your first step is labeling everything. Knowing you use size to predict price, you figure size must be the *x* variable and price must be the *y* variable. You then label the means $\bar{x} = 1,993$ (square feet) and $\bar{y} = 219.1$ (in thousands) respectively; the standard deviations are labeled $s_x = 349$ (square feet) and $s_y = 60.1$ (in thousands), respectively, and the correlation is labeled $r = 0.90$. The sample size is $n = 100$. Now you can plug your numbers into the right formulas. (See Chapter 18 regarding correlation and regression.)

When you know you have to work with a regression line and that formulas are involved, having all the given information organized and labeled, ready to go, is very comforting. It's one less thing to think about. (The problem in this particular example is solved in the section "Make the Connection and Solve the Problem.") If that example doesn't convince you, here are six more reasons to label what you are given in a problem:

- ✓ **Labeling allows you to check your work more easily.** When you go back to check your work (as I advise in the section "Do the Math — Twice"), you'll quickly see what you were thinking when you did the problem the first time.
- ✓ **Your professor will be impressed.** He will see your labels and realize you at least know what the given information stands for. That way if your calculations go haywire, you still have a chance for partial credit.
- ✓ **Labeling saves time.** I know that writing down more information seems like a strange way to save time, but by labeling all the items, you can pull out the info you need in a flash.

For example, suppose you need to do a 95% confidence interval for the population mean (using what you know from Chapter 13) and you're told that the sample mean is 60, the population standard deviation is 10, and the sample size is 200. You know the formula has to involve \bar{x} , σ , and n , and you see one that does:

$$\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$$

Because you've already labeled everything, you just grab what you need, put it into the formula, throw in a z^* -value of 1.96 (the critical value corresponding to a 95% confidence level), and crunch it out to get the answer:

$$60 \pm 1.96 \frac{10}{\sqrt{200}} = 60 \pm 1.39$$

- ✓ **Labels keep your mind organized.** You are less likely to get buried in calculations and forget what you're doing if your work involves symbols and not just numbers. By sorting out the information you are given, you're less likely to resort to reading the problem over and over again, raising your anxiety level each time.
- ✓ **You use the labels to figure out which formula or technique you need to use to solve a problem.** For example, if you think you need a hypothesis test but no claim is made about the population mean, hold up. You may need a confidence interval instead; this realization saves you precious time because you won't be spinning your wheels in the wrong direction. Labels help you quickly narrow down your options.
- ✓ **Labeling helps you resist the urge to just write down numbers and push them around on the paper.** More often than not, number-pushing leads to wrong answers and less (if any) partial credit if your answer is wrong. Your professor may not be able to follow you, or just doesn't want to spend all that time trying to

figure it out (sorry to say, but this happens sometimes).



Labeling saves you anxiety, time, and points when you take your exam. But in order to be successful on exam day, you need to start this practice early on, while the problems are easy to do. Don't expect to suddenly be able to sort out the information on exam day if you never did it before; it's not gonna happen. Make it your habit right away and you won't freak out when you see a new problem. You'll at least be able to break it down into smaller chunks, which always helps.

Draw a Picture

You've heard the expression "A picture is worth a thousand words." As a statistics professor, I say, "A picture is worth a thousand points (or at least half the points on a given problem)." When the given information and/or the question being asked can be expressed in a picture form, you should do it. Here's why:

- ✓ **A picture can help you see what's going on in the problem.** For example, if you know exam scores have a normal distribution with mean 75 and standard deviation 5 (see Chapter 9 for more about normal distribution), you draw a bell-shaped curve, marking off the mean in the center and three standard deviations on each side. You can now visualize the scenario you're dealing with.
- ✓ **You can use the drawing to help figure out what you are trying to find.** For example, if you need to know the probability that Bob scored more than 70 points on the exam, you shade in the area to the right of 70 on your drawing, and you're on your way.
- ✓ **Your professor knows that you understand the basics of the problem, increasing your chance for partial credit.** On the other hand, someone who got the problem wrong doesn't get much sympathy if the professor knows drawing a simple picture would have avoided the whole problem.
- ✓ **Students who draw pictures tend to get more problems correct than students who don't.** Without a picture you can easily lose track of what's needed, and make mistakes like finding $P(X < 70)$ instead of $P(X > 70)$, for example. Also, checking for and spotting errors before you turn in your exam is easier if you have a picture to look at.



Drawing a picture may seem like a waste of valuable time on an exam, but it's actually a time-saver because it gets you going in the right direction, keeps you focused throughout the problem, and helps ensure you answer the right question. Drawing a picture can also help you analyze your final numerical answer and either

confirm you've got it right, or quickly a spot and fix an error and save yourself some points. (Be sure to draw pictures while studying so they come naturally during an exam.)

Make the Connection and Solve the Problem



When you've figured out what the problem is asking, you have everything labeled, and you have your pictures drawn, it's time to solve the problem. After doing the prep work, nine times out of ten you'll remember a technique you learned from class, a formula that contains the items you've labeled, and/or an example you worked through. Use or remember your if-then-how chart and you'll be on your way. (See "Make an 'If-Then-How' Chart" if you need more info.)



Breaking down a problem means having less to think about at each step, and in a stressful exam situation where you may forget your own name, that's a real plus! (This strategy reminds me of the saying, "How do you eat an elephant? One bite at a time.")

In the example of using size of a home to predict its price (see the earlier section "Label What You're Given"), you know the mean and standard deviation of size, the mean and standard deviation of price, and the correlation between them; and you've labeled them all. The question asks you to find the equation of the best-fitting regression line to predict price based on size of the home; you know that means find the equation $y = a + bx$ where x = size (square feet) and y = price (thousands of dollars), b is the slope of the regression line, and a is the y -intercept. (Flip to Chapter 18 for more about this formula.)

Now you recognize what to do — you have to find a and b . You remember (or can find) that those formulas are $b = r \frac{s_y}{s_x}$ and $a = \bar{y} - b\bar{x}$. Grab the numbers you've labeled ($\bar{x} = 1,993$; $s_x = 349$; $\bar{y} = 219.1$; $s_y = 60.1$; and $r = 0.90$), put them into the formulas, and solve (sounds like a commercial for a frozen dinner doesn't it?). You find the slope is $b = 0.90 \frac{60.1}{349} = 0.155$ and the y -intercept is $a = 219.1 - 0.155(1,993) = -89.82$, so the equation of the best-fitting regression line is $y = -89.82 + 0.155x$. (See Chapter 18 for the details of regression.)

Do the Math — Twice

I can still remember some of the struggles I had way back in high school algebra. For the longest time 3 times 2 was equal to 5 for me; this mistake (and others like it) caused me to miss a handful of points on every exam and homework assignment, and I just could

not get past it. One day I decided I'd had enough of losing points here and there for silly errors, and I did something about it. From that day on, I wrote out all of my work, step by step, and resisted the urge to do steps in my head. When I got my final answer, instead of moving on, I went back and checked every step, and I did so with the mind-set that a mistake had probably slipped in somewhere and it was my job to find it before anyone else did.

This approach forced me to look at each step with fresh eyes, as if I were grading someone else's paper. I caught more mistakes because I never skipped over a step without bothering to check it. I finally stopped thinking 3 times 2 was 5 because I caught myself in the act enough times. My exam grades went up, just because I started checking things more carefully. It reminds me of the carpenter's saying, "Measure twice, cut once." They waste a lot less wood that way.



Every time you find and fix a mistake before you turn in your exam, you're getting a handful of points back for yourself. Find your errors before your professor does, and you'll be amazed how those points add up. However, remember that time is not unlimited on an exam, so try to get the problems right the first time. Labeling everything, drawing pictures, writing down formulas, and showing all your work will definitely help!

Analyze Your Answers

A very prominent statistician I know has a framed piece of paper on his office wall. It's a page of an exam he took way back when he was a student. It's got a big red circle around one of his answers, which happens to be the number 2. Why was writing the number 2 for an answer such a problem? Because the question asked him to find a probability, and probabilities are always between 0 and 1. As a result, he didn't get any points for that problem, not even partial credit. In fact, I'll bet his professor wanted to give him negative points for making such a mistake. (They really don't like it when you totally miss the boat.)



Always take the time to check your final answer to see if it makes sense. A negative standard deviation, a probability more than 1, or a correlation of -121.23 is not going to go over well with your professor, and it will not be treated like a simple math error. It will be treated as a fundamental error in not knowing (or perhaps caring) what the result should look like.



If you know an answer you got can't possibly be right, but you cannot for the life of you figure out where you went wrong, don't waste any more time on it. Just write a

note in the margin that says you know your answer can't be right but you can't figure out your error. This helps separate you from the regular Joe who found a probability of 10,524.31 (yes, I've seen it) and merrily moved on.

By the way, you may be wondering why this world-class statistician still keeps this exam page framed on his office wall. He says it's to keep him humble. Learn from his example and never move on to the next problem without stepping back and saying "does this answer even make sense?"

Appendix

Tables for Reference

This appendix includes tables for finding probabilities and/or critical values for the three distributions used in this book: the Z-distribution (standard normal), the *t*-distribution, and the binomial distribution.

The Z-Table

Table A-1 shows less-than-or-equal-to probabilities for the Z-distribution; that is, $p(Z \leq z)$ for a given z -value. (See Chapter 9 for calculating z -values for a normal distribution; see Chapter 11 for calculating z -values for a sampling distribution.) To use Table A-1, do the following:

1. Determine the z -value for your particular problem.

The z -value should have one leading digit before the decimal point (positive, negative, or zero) and two digits after the decimal point; for example $z = 1.28$, -2.69 , or 0.13 .

2. Find the row of the table corresponding to the leading digit and first digit after the decimal point.

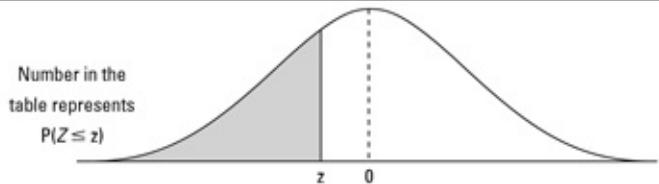
For example, if your z -value is 1.28 , look in the “ 1.2 ” row; if $z = -1.28$, look in the “ -1.2 ” row.

3. Find the column corresponding to the second digit after the decimal point.

For example, if your z -value is 1.28 or -1.28 , look in the “ $.08$ ” column.

4. Intersect the row and column from Steps 2 and 3.

This number is the probability that Z is less than or equal to your z -value. In other words, you’ve found $p(Z \leq z)$. For example, if $z = 1.28$, you see $p(Z \leq 1.28) = 0.8997$. For $z = -1.28$, you see $p(Z \leq -1.28) = 0.1003$.

Table A-1**The Z-Table**

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
-3.6	.0002	.0002	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001
-3.5	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002
-3.4	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0002	.0002
-3.3	.0005	.0005	.0005	.0004	.0004	.0004	.0004	.0004	.0003	.0003
-3.2	.0007	.0007	.0006	.0006	.0006	.0006	.0006	.0005	.0005	.0005
-3.1	.0010	.0009	.0009	.0009	.0008	.0008	.0008	.0008	.0007	.0007
-3.0	.0013	.0013	.0013	.0012	.0012	.0011	.0011	.0011	.0010	.0010
-2.9	.0019	.0018	.0018	.0017	.0016	.0016	.0015	.0015	.0014	.0014
-2.8	.0026	.0025	.0024	.0023	.0023	.0022	.0021	.0021	.0020	.0019
-2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
-2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
-2.5	.0062	.0060	.0059	.0057	.0055	.0054	.0052	.0051	.0049	.0048
-2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064
-2.3	.0107	.0104	.0102	.0099	.0096	.0094	.0091	.0089	.0087	.0084
-2.2	.0139	.0136	.0132	.0129	.0125	.0122	.0119	.0116	.0113	.0110
-2.1	.0179	.0174	.0170	.0166	.0162	.0158	.0154	.0150	.0146	.0143
-2.0	.0228	.0222	.0217	.0212	.0207	.0202	.0197	.0192	.0188	.0183
-1.9	.0287	.0281	.0274	.0268	.0262	.0256	.0250	.0244	.0239	.0233
-1.8	.0359	.0351	.0344	.0336	.0329	.0322	.0314	.0307	.0301	.0294
-1.7	.0446	.0436	.0427	.0418	.0409	.0401	.0392	.0384	.0375	.0367
-1.6	.0548	.0537	.0526	.0516	.0505	.0495	.0485	.0475	.0465	.0455
-1.5	.0668	.0655	.0643	.0630	.0618	.0606	.0594	.0582	.0571	.0559
-1.4	.0808	.0793	.0778	.0764	.0749	.0735	.0721	.0708	.0694	.0681
-1.3	.0968	.0951	.0934	.0918	.0901	.0885	.0869	.0853	.0838	.0823
-1.2	.1151	.1131	.1112	.1093	.1075	.1056	.1038	.1020	.1003	.0985
-1.1	.1357	.1335	.1314	.1292	.1271	.1251	.1230	.1210	.1190	.1170
-1.0	.1587	.1562	.1539	.1515	.1492	.1469	.1446	.1423	.1401	.1379
-0.9	.1841	.1814	.1788	.1762	.1736	.1711	.1685	.1660	.1635	.1611
-0.8	.2119	.2090	.2061	.2033	.2005	.1977	.1949	.1922	.1894	.1867
-0.7	.2420	.2389	.2358	.2327	.2296	.2266	.2236	.2206	.2177	.2148
-0.6	.2743	.2709	.2676	.2643	.2611	.2578	.2546	.2514	.2483	.2451
-0.5	.3085	.3050	.3015	.2981	.2946	.2912	.2877	.2843	.2810	.2776
-0.4	.3446	.3409	.3372	.3336	.3300	.3264	.3228	.3192	.3156	.3121
-0.3	.3821	.3783	.3745	.3707	.3669	.3632	.3594	.3557	.3520	.3483
-0.2	.4207	.4168	.4129	.4090	.4052	.4013	.3974	.3936	.3897	.3859
-0.1	.4602	.4562	.4522	.4483	.4443	.4404	.4364	.4325	.4286	.4247
-0.0	.5000	.4960	.4920	.4880	.4840	.4801	.4761	.4721	.4681	.4641

Number in the table represents $P(Z \leq z)$

<i>z</i>	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
0.1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
0.2	.5793	.5832	.5871	.5910	.5948	.5987	.6026	.6064	.6103	.6141
0.3	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517
0.4	.6554	.6591	.6628	.6664	.6700	.6736	.6772	.6808	.6844	.6879
0.5	.6915	.6950	.6985	.7019	.7054	.7088	.7123	.7157	.7190	.7224
0.6	.7257	.7291	.7324	.7357	.7389	.7422	.7454	.7486	.7517	.7549
0.7	.7580	.7611	.7642	.7673	.7704	.7734	.7764	.7794	.7823	.7852
0.8	.7881	.7910	.7939	.7967	.7995	.8023	.8051	.8078	.8106	.8133
0.9	.8159	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.8461	.8485	.8508	.8531	.8554	.8577	.8599	.8621
1.1	.8643	.8665	.8686	.8708	.8729	.8749	.8770	.8790	.8810	.8830
1.2	.8849	.8869	.8888	.8907	.8925	.8944	.8962	.8980	.8997	.9015
1.3	.9032	.9049	.9066	.9082	.9099	.9115	.9131	.9147	.9162	.9177
1.4	.9192	.9207	.9222	.9236	.9251	.9265	.9279	.9292	.9306	.9319
1.5	.9332	.9345	.9357	.9370	.9382	.9394	.9406	.9418	.9429	.9441
1.6	.9452	.9463	.9474	.9484	.9495	.9505	.9515	.9525	.9535	.9545
1.7	.9554	.9564	.9573	.9582	.9591	.9599	.9608	.9616	.9625	.9633
1.8	.9641	.9649	.9656	.9664	.9671	.9678	.9686	.9693	.9699	.9706
1.9	.9713	.9719	.9726	.9732	.9738	.9744	.9750	.9756	.9761	.9767
2.0	.9772	.9778	.9783	.9788	.9793	.9798	.9803	.9808	.9812	.9817
2.1	.9821	.9826	.9830	.9834	.9838	.9842	.9846	.9850	.9854	.9857
2.2	.9861	.9864	.9868	.9871	.9875	.9878	.9881	.9884	.9887	.9890
2.3	.9893	.9896	.9898	.9901	.9904	.9906	.9909	.9911	.9913	.9916
2.4	.9918	.9920	.9922	.9925	.9927	.9929	.9931	.9932	.9934	.9936
2.5	.9938	.9940	.9941	.9943	.9945	.9946	.9948	.9949	.9951	.9952
2.6	.9953	.9955	.9956	.9957	.9959	.9960	.9961	.9962	.9963	.9964
2.7	.9965	.9966	.9967	.9968	.9969	.9970	.9971	.9972	.9973	.9974
2.8	.9974	.9975	.9976	.9977	.9977	.9978	.9979	.9979	.9980	.9981
2.9	.9981	.9982	.9982	.9983	.9984	.9984	.9985	.9985	.9986	.9986
3.0	.9987	.9987	.9987	.9988	.9988	.9989	.9989	.9989	.9990	.9990
3.1	.9990	.9991	.9991	.9991	.9992	.9992	.9992	.9992	.9993	.9993
3.2	.9993	.9993	.9994	.9994	.9994	.9994	.9994	.9995	.9995	.9995
3.3	.9995	.9995	.9995	.9996	.9996	.9996	.9996	.9996	.9996	.9997
3.4	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9998
3.5	.9998	.9998	.9998	.9998	.9998	.9998	.9998	.9998	.9998	.9998
3.6	.9998	.9998	.9999	.9999	.9999	.9999	.9999	.9999	.9999	.9999

The t-Table

Table A-2 shows right-tail probabilities for selected *t*-distributions (see Chapter 10 for more on the *t*-distribution).

Follow these steps to use Table A-2 to find right-tail probabilities and *p*-values for hypothesis tests involving *t* (see Chapter 15):

1. Find the *t*-value for which you want the right-tail probability (call it *t*), and find the sample size (for example, *n*).

2. Find the row corresponding to the degrees of freedom (df) for your problem (for example, $n - 1$). Go across that row to find the two t -values between which your t falls.

For example, if your t is 1.60 and your n is 7, you look in the row for $df = 7 - 1 = 6$. Across that row you find your t lies between t -values 1.44 and 1.94.

3. Go to the top of the columns containing the two t -values from Step 2.

The right-tail (greater-than) probability for your t -value is somewhere between the two values at the top of these columns. For example, your $t = 1.60$ is between t -values 1.44 and 1.94 ($df = 6$); so the right tail probability for your t is between 0.10 (column heading for $t = 1.44$); and 0.05 (column heading for $t = 1.94$).



The row near the bottom with Z in the df column gives right-tail (greater-than) probabilities from the Z -distribution (Chapter 10 shows Z 's relationship with t).

Use Table A-2 to find t^* -values (critical values) for a confidence interval involving t (see Chapter 13):

- 1. Determine the confidence level you need (as a percentage).**
- 2. Determine the sample size (for example, n).**
- 3. Look at the bottom row of the table where the percentages are shown. Find your % confidence level there.**

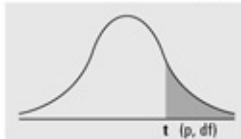
4. Intersect this column with the row representing your degrees of freedom (df).

This is the t -value you need for your confidence interval.

For example, a 95% confidence interval with $df=6$ has $t^*=2.45$. (Find 95% on the last line and go up to row 6.)

Table A-2**The *t*-Table**

Numbers in each row of the table are values on a *t*-distribution with (*df*) degrees of freedom for selected right-tail (greater-than) probabilities (*p*).



df/p	0.40	0.25	0.10	0.05	0.025	0.01	0.005	0.0005
1	0.324920	1.000000	3.077684	6.313752	12.70620	31.82052	63.65674	636.6192
2	0.288675	0.816497	1.885618	2.919986	4.30265	6.96456	9.92484	31.5991
3	0.276671	0.764892	1.637744	2.353363	3.18245	4.54070	5.84091	12.9240
4	0.270722	0.740697	1.533206	2.131847	2.77645	3.74695	4.60409	8.6103
5	0.267181	0.726687	1.475884	2.015048	2.57058	3.36493	4.03214	6.8688
6	0.264835	0.717558	1.439756	1.943180	2.44691	3.14267	3.70743	5.9588
7	0.263167	0.711142	1.414924	1.894579	2.36462	2.99795	3.49948	5.4079
8	0.261921	0.706387	1.396815	1.859548	2.30600	2.89646	3.35539	5.0413
9	0.260955	0.702722	1.383029	1.833113	2.26216	2.82144	3.24984	4.7809
10	0.260185	0.699812	1.372184	1.812461	2.22814	2.76377	3.16927	4.5869
11	0.259556	0.697445	1.363430	1.795885	2.20099	2.71808	3.10581	4.4370
12	0.259033	0.695483	1.356217	1.782288	2.17881	2.68100	3.05454	43178
13	0.258591	0.693829	1.350171	1.770933	2.16037	2.65031	3.01228	4.2208
14	0.258213	0.692417	1.345030	1.761310	2.14479	2.62449	2.97684	4.1405
15	0.257885	0.691197	1.340606	1.753050	2.13145	2.60248	2.94671	4.0728
16	0.257599	0.690132	1.336757	1.745884	2.11991	2.58349	2.92078	4.0150
17	0.257347	0.689195	1.333379	1.739607	2.10982	2.56693	2.89823	3.9651
18	0.257123	0.688364	1.330391	1.734064	2.10092	2.55238	2.87844	3.9216
19	0.256923	0.687621	1.327728	1.729133	2.09302	2.53948	2.86093	3.8834
20	0.256743	0.686954	1.325341	1.724718	2.08596	2.52798	2.84534	3.8495
21	0.256580	0.686352	1.323188	1.720743	2.07961	2.51765	2.83136	3.8193
22	0.256432	0.685805	1.321237	1.717144	2.07387	2.50832	2.81876	3.7921
23	0.256297	0.685306	1.319460	1.713872	2.06866	2.49987	2.80734	3.7676
24	0.256173	0.684850	1.317836	1.710882	2.06390	2.49216	2.79694	3.7454
25	0.256060	0.684430	1.316345	1.708141	2.05954	2.48511	2.78744	3.7251
26	0.255955	0.684043	1.314972	1.705618	2.05553	2.47863	2.77871	3.7066
27	0.255858	0.683685	1.313703	1.703288	2.05183	2.47266	2.77068	3.6896
28	0.255768	0.683353	1.312527	1.701131	2.04841	2.46714	2.76326	3.6739
29	0.255684	0.683044	1.311434	1.699127	2.04523	2.46202	2.75639	3.6594
30	0.255605	0.682756	1.310415	1.697261	2.04227	2.45726	2.75000	3.6460
z	0.253347	0.674490	1.281552	1.644854	1.95996	2.32635	2.57583	3.2905
CI	—	—	80%	90%	95%	98%	99%	99.9%

The Binomial Table

Table A-3 shows probabilities for the binomial distribution (see Chapter 8).

To use Table A-3, do the following:

1. Find these three numbers for your particular problem:

- The sample size, n
- The probability of success, p
- The x -value for which you want $p(X = x)$

2. Find the section of Table A-3 that's devoted to your n .

3. Look at the row for your x -value and the column for your p .

4. Intersect that row and column. You have found $p(X = x)$.

5. To get the probability of being less than, greater than, greater than or equal to, less than or equal to, or between two values of X , you add the appropriate values of Table A-3 using the steps found in Chapter 8.

For example, if $n=10$, $p=0.6$, and you want $p(X=9)$, go to the $n=10$ section, the $x=9$ row, and the $p=0.6$ column to find 0.04.

Table A-3

The Binomial Table

Numbers in the table represent $p(X=x)$ for a binomial distribution with n trials and probability of success p .

Binomial probabilities: $\binom{n}{x} p^x(1-p)^{n-x}$		p										
n	x	0.1	0.2	0.25	0.3	0.4	0.5	0.6	0.7	0.75	0.8	0.9
1	0	0.900	0.800	0.750	0.700	0.600	0.500	0.400	0.300	0.250	0.200	0.100
	1	0.100	0.200	0.250	0.300	0.400	0.500	0.600	0.700	0.750	0.800	0.900
2	0	0.810	0.640	0.563	0.490	0.360	0.250	0.160	0.090	0.063	0.040	0.010
	1	0.180	0.320	0.375	0.420	0.480	0.500	0.480	0.420	0.375	0.320	0.180
	2	0.010	0.040	0.063	0.090	0.160	0.250	0.360	0.490	0.563	0.640	0.810
3	0	0.729	0.512	0.422	0.343	0.216	0.125	0.064	0.027	0.016	0.008	0.001
	1	0.243	0.384	0.422	0.441	0.432	0.375	0.288	0.189	0.141	0.096	0.027
	2	0.027	0.096	0.141	0.189	0.288	0.375	0.432	0.441	0.422	0.384	0.243
	3	0.001	0.008	0.016	0.027	0.064	0.125	0.216	0.343	0.422	0.512	0.729
4	0	0.656	0.410	0.316	0.240	0.130	0.063	0.026	0.008	0.004	0.002	0.000
	1	0.292	0.410	0.422	0.412	0.346	0.250	0.154	0.076	0.047	0.026	0.004
	2	0.049	0.154	0.211	0.265	0.346	0.375	0.346	0.265	0.211	0.154	0.049
	3	0.004	0.026	0.047	0.076	0.154	0.250	0.346	0.412	0.422	0.410	0.292
	4	0.000	0.002	0.004	0.008	0.026	0.063	0.130	0.240	0.316	0.410	0.656
5	0	0.590	0.328	0.237	0.168	0.078	0.031	0.010	0.002	0.001	0.000	0.000
	1	0.328	0.410	0.396	0.360	0.259	0.156	0.077	0.028	0.015	0.006	0.000
	2	0.073	0.205	0.264	0.309	0.346	0.312	0.230	0.132	0.088	0.051	0.008
	3	0.008	0.051	0.088	0.132	0.230	0.312	0.346	0.309	0.264	0.205	0.073
	4	0.000	0.006	0.015	0.028	0.077	0.156	0.259	0.360	0.396	0.410	0.328
	5	0.000	0.000	0.001	0.002	0.010	0.031	0.078	0.168	0.237	0.328	0.590
6	0	0.531	0.262	0.178	0.118	0.047	0.016	0.004	0.001	0.000	0.000	0.000
	1	0.354	0.393	0.356	0.303	0.187	0.094	0.037	0.010	0.004	0.002	0.000
	2	0.098	0.246	0.297	0.324	0.311	0.234	0.138	0.060	0.033	0.015	0.001
	3	0.015	0.082	0.132	0.185	0.276	0.313	0.276	0.185	0.132	0.082	0.015
	4	0.001	0.015	0.033	0.060	0.138	0.234	0.311	0.324	0.297	0.246	0.098
	5	0.000	0.002	0.004	0.010	0.037	0.094	0.187	0.303	0.356	0.393	0.354
	6	0.000	0.000	0.000	0.001	0.004	0.016	0.047	0.118	0.178	0.262	0.531
7	0	0.478	0.210	0.133	0.082	0.028	0.008	0.002	0.000	0.000	0.000	0.000
	1	0.372	0.367	0.311	0.247	0.131	0.055	0.017	0.004	0.001	0.000	0.000
	2	0.124	0.275	0.311	0.318	0.261	0.164	0.077	0.025	0.012	0.004	0.000
	3	0.023	0.115	0.173	0.227	0.290	0.273	0.194	0.097	0.058	0.029	0.003
	4	0.003	0.029	0.058	0.097	0.194	0.273	0.290	0.227	0.173	0.115	0.023
	5	0.000	0.004	0.012	0.025	0.077	0.164	0.261	0.318	0.311	0.275	0.124
	6	0.000	0.000	0.001	0.004	0.017	0.055	0.131	0.247	0.311	0.367	0.372
	7	0.000	0.000	0.000	0.000	0.002	0.008	0.028	0.082	0.133	0.210	0.478

Table A-3

Numbers in the table represent $p(X=x)$ for a binomial distribution with n trials and probability of success p .

		p										
		0.1	0.2	0.25	0.3	0.4	0.5	0.6	0.7	0.75	0.8	0.9
Binomial probabilities: $\binom{n}{x} p^x (1-p)^{n-x}$		<i>n</i>	<i>x</i>									
$\binom{n}{x} p^x (1-p)^{n-x}$		8	0	0.430	0.168	0.100	0.058	0.017	0.004	0.001	0.000	0.000
		8	1	0.383	0.336	0.267	0.198	0.090	0.031	0.008	0.001	0.000
		8	2	0.149	0.294	0.311	0.296	0.209	0.109	0.041	0.010	0.004
		8	3	0.033	0.147	0.208	0.254	0.279	0.219	0.124	0.047	0.023
		8	4	0.005	0.046	0.087	0.136	0.232	0.273	0.232	0.136	0.087
		8	5	0.000	0.009	0.023	0.047	0.124	0.219	0.279	0.254	0.208
		8	6	0.000	0.001	0.004	0.010	0.041	0.109	0.209	0.296	0.311
		8	7	0.000	0.000	0.000	0.001	0.008	0.031	0.090	0.198	0.267
		8	8	0.000	0.000	0.000	0.000	0.001	0.004	0.017	0.058	0.100
$\binom{n}{x} p^x (1-p)^{n-x}$		9	0	0.387	0.134	0.075	0.040	0.010	0.002	0.000	0.000	0.000
		9	1	0.387	0.302	0.225	0.156	0.060	0.018	0.004	0.000	0.000
		9	2	0.172	0.302	0.300	0.267	0.161	0.070	0.021	0.004	0.001
		9	3	0.045	0.176	0.234	0.267	0.251	0.164	0.074	0.021	0.009
		9	4	0.007	0.066	0.117	0.172	0.251	0.246	0.167	0.074	0.039
		9	5	0.001	0.017	0.039	0.074	0.167	0.246	0.251	0.172	0.117
		9	6	0.000	0.003	0.009	0.021	0.074	0.164	0.251	0.267	0.234
		9	7	0.000	0.000	0.001	0.004	0.021	0.070	0.161	0.267	0.300
		9	8	0.000	0.000	0.000	0.000	0.004	0.018	0.060	0.156	0.225
		9	9	0.000	0.000	0.000	0.000	0.000	0.002	0.010	0.040	0.075
$\binom{n}{x} p^x (1-p)^{n-x}$		10	0	0.349	0.107	0.056	0.028	0.006	0.001	0.000	0.000	0.000
		10	1	0.387	0.268	0.188	0.121	0.040	0.010	0.002	0.000	0.000
		10	2	0.194	0.302	0.282	0.233	0.121	0.044	0.011	0.001	0.000
		10	3	0.057	0.201	0.250	0.267	0.215	0.117	0.042	0.009	0.003
		10	4	0.011	0.088	0.146	0.200	0.251	0.205	0.111	0.037	0.016
		10	5	0.001	0.026	0.058	0.103	0.201	0.246	0.201	0.103	0.058
		10	6	0.000	0.006	0.016	0.037	0.111	0.205	0.251	0.200	0.146
		10	7	0.000	0.001	0.003	0.009	0.042	0.117	0.215	0.267	0.250
		10	8	0.000	0.000	0.000	0.001	0.011	0.044	0.121	0.233	0.282
		10	9	0.000	0.000	0.000	0.000	0.002	0.010	0.040	0.121	0.188
		10	10	0.000	0.000	0.000	0.000	0.000	0.001	0.006	0.028	0.056
		10	11	0.000	0.000	0.000	0.000	0.000	0.000	0.006	0.017	0.049
$\binom{n}{x} p^x (1-p)^{n-x}$		11	0	0.314	0.086	0.042	0.020	0.004	0.000	0.000	0.000	0.000
		11	1	0.384	0.236	0.155	0.093	0.027	0.005	0.001	0.000	0.000
		11	2	0.213	0.295	0.258	0.200	0.089	0.027	0.005	0.001	0.000
		11	3	0.071	0.221	0.258	0.257	0.177	0.081	0.023	0.004	0.001
		11	4	0.016	0.111	0.172	0.220	0.236	0.161	0.070	0.017	0.006
		11	5	0.002	0.039	0.080	0.132	0.221	0.226	0.147	0.057	0.027
		11	6	0.000	0.010	0.027	0.057	0.147	0.226	0.221	0.132	0.080
		11	7	0.000	0.002	0.006	0.017	0.070	0.161	0.236	0.220	0.172
		11	8	0.000	0.000	0.001	0.004	0.023	0.081	0.177	0.257	0.258
		11	9	0.000	0.000	0.000	0.001	0.005	0.027	0.089	0.200	0.258
		11	10	0.000	0.000	0.000	0.000	0.001	0.005	0.027	0.093	0.155
		11	11	0.000	0.000	0.000	0.000	0.000	0.004	0.020	0.042	0.086

Numbers in the table represent $p(X=x)$ for a binomial distribution with n trials and probability of success p .

Binomial probabilities:

$$\binom{n}{x} p^x (1-p)^{n-x}$$

		p										
n	x	0.1	0.2	0.25	0.3	0.4	0.5	0.6	0.7	0.75	0.8	0.9
12	0	0.282	0.069	0.032	0.014	0.002	0.000	0.000	0.000	0.000	0.000	0.000
	1	0.377	0.206	0.127	0.071	0.017	0.003	0.000	0.000	0.000	0.000	0.000
	2	0.230	0.283	0.232	0.168	0.064	0.016	0.002	0.000	0.000	0.000	0.000
	3	0.085	0.236	0.258	0.240	0.142	0.054	0.012	0.001	0.000	0.000	0.000
	4	0.021	0.133	0.194	0.231	0.213	0.121	0.042	0.008	0.002	0.001	0.000
	5	0.004	0.053	0.103	0.158	0.227	0.193	0.101	0.029	0.011	0.003	0.000
	6	0.000	0.016	0.040	0.079	0.177	0.226	0.177	0.079	0.040	0.016	0.000
	7	0.000	0.003	0.011	0.029	0.101	0.193	0.227	0.158	0.103	0.053	0.004
	8	0.000	0.001	0.002	0.008	0.042	0.121	0.213	0.231	0.194	0.133	0.021
	9	0.000	0.000	0.000	0.001	0.012	0.054	0.142	0.240	0.258	0.236	0.085
	10	0.000	0.000	0.000	0.000	0.002	0.016	0.064	0.168	0.232	0.283	0.230
	11	0.000	0.000	0.000	0.000	0.000	0.003	0.017	0.071	0.127	0.206	0.377
	12	0.000	0.000	0.000	0.000	0.000	0.000	0.002	0.014	0.032	0.069	0.282
13	0	0.254	0.055	0.024	0.010	0.001	0.000	0.000	0.000	0.000	0.000	0.000
	1	0.367	0.179	0.103	0.054	0.011	0.002	0.000	0.000	0.000	0.000	0.000
	2	0.245	0.268	0.206	0.139	0.045	0.010	0.001	0.000	0.000	0.000	0.000
	3	0.100	0.246	0.252	0.218	0.111	0.035	0.006	0.001	0.000	0.000	0.000
	4	0.028	0.154	0.210	0.234	0.184	0.087	0.024	0.003	0.001	0.000	0.000
	5	0.006	0.069	0.126	0.180	0.221	0.157	0.066	0.014	0.005	0.001	0.000
	6	0.001	0.023	0.056	0.103	0.197	0.209	0.131	0.044	0.019	0.006	0.000
	7	0.000	0.006	0.019	0.044	0.131	0.209	0.197	0.103	0.056	0.023	0.001
	8	0.000	0.001	0.005	0.014	0.066	0.157	0.221	0.180	0.126	0.069	0.006
	9	0.000	0.000	0.001	0.003	0.024	0.087	0.184	0.234	0.210	0.154	0.028
	10	0.000	0.000	0.000	0.001	0.006	0.035	0.111	0.218	0.252	0.246	0.100
	11	0.000	0.000	0.000	0.000	0.001	0.010	0.045	0.139	0.206	0.268	0.245
	12	0.000	0.000	0.000	0.000	0.000	0.002	0.011	0.054	0.103	0.179	0.367
	13	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.010	0.024	0.055	0.254
14	0	0.229	0.044	0.018	0.007	0.001	0.000	0.000	0.000	0.000	0.000	0.000
	1	0.356	0.154	0.083	0.041	0.007	0.001	0.000	0.000	0.000	0.000	0.000
	2	0.257	0.250	0.180	0.113	0.032	0.006	0.001	0.000	0.000	0.000	0.000
	3	0.114	0.250	0.240	0.194	0.085	0.022	0.003	0.000	0.000	0.000	0.000
	4	0.035	0.172	0.220	0.229	0.155	0.061	0.014	0.001	0.000	0.000	0.000
	5	0.008	0.086	0.147	0.196	0.207	0.122	0.041	0.007	0.002	0.000	0.000
	6	0.001	0.032	0.073	0.126	0.207	0.183	0.092	0.023	0.008	0.002	0.000
	7	0.000	0.009	0.028	0.062	0.157	0.209	0.157	0.062	0.028	0.009	0.000
	8	0.000	0.002	0.008	0.023	0.092	0.183	0.207	0.126	0.073	0.032	0.001
	9	0.000	0.000	0.002	0.007	0.041	0.122	0.207	0.196	0.147	0.086	0.008
	10	0.000	0.000	0.000	0.001	0.014	0.061	0.155	0.229	0.220	0.172	0.035
	11	0.000	0.000	0.000	0.000	0.003	0.022	0.085	0.194	0.240	0.250	0.114
	12	0.000	0.000	0.000	0.000	0.001	0.006	0.032	0.113	0.180	0.250	0.257
	13	0.000	0.000	0.000	0.000	0.000	0.001	0.007	0.041	0.083	0.154	0.356
	14	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.007	0.018	0.044	0.229

Table A-3

Numbers in the table represent $p(X=x)$ for a binomial distribution with n trials and probability of success p .

		p										
		0.1	0.2	0.25	0.3	0.4	0.5	0.6	0.7	0.75	0.8	0.9
Binomial probabilities:		$\binom{n}{x} p^x (1-p)^{n-x}$										
n	x											
15	0	0.206	0.035	0.013	0.005	0.000	0.000	0.000	0.000	0.000	0.000	
	1	0.343	0.132	0.067	0.031	0.005	0.000	0.000	0.000	0.000	0.000	
	2	0.267	0.231	0.156	0.092	0.022	0.003	0.000	0.000	0.000	0.000	
	3	0.129	0.250	0.225	0.170	0.063	0.014	0.002	0.000	0.000	0.000	
	4	0.043	0.188	0.225	0.219	0.127	0.042	0.007	0.001	0.000	0.000	
	5	0.010	0.103	0.165	0.206	0.186	0.092	0.024	0.003	0.001	0.000	
	6	0.002	0.043	0.092	0.147	0.207	0.153	0.061	0.012	0.003	0.001	
	7	0.000	0.014	0.039	0.081	0.177	0.196	0.118	0.035	0.013	0.003	
	8	0.000	0.003	0.013	0.035	0.118	0.196	0.177	0.081	0.039	0.014	
	9	0.000	0.001	0.003	0.012	0.061	0.153	0.207	0.147	0.092	0.043	
	10	0.000	0.000	0.001	0.003	0.024	0.092	0.186	0.206	0.165	0.103	
	11	0.000	0.000	0.000	0.001	0.007	0.042	0.127	0.219	0.225	0.188	
	12	0.000	0.000	0.000	0.000	0.002	0.014	0.063	0.170	0.225	0.250	
	13	0.000	0.000	0.000	0.000	0.000	0.003	0.022	0.092	0.156	0.231	
	14	0.000	0.000	0.000	0.000	0.000	0.000	0.005	0.031	0.067	0.132	
	15	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.005	0.013	0.035	
20	0	0.122	0.012	0.003	0.001	0.000	0.000	0.000	0.000	0.000	0.000	
	1	0.270	0.058	0.021	0.007	0.000	0.000	0.000	0.000	0.000	0.000	
	2	0.285	0.137	0.067	0.028	0.003	0.000	0.000	0.000	0.000	0.000	
	3	0.190	0.205	0.134	0.072	0.012	0.001	0.000	0.000	0.000	0.000	
	4	0.090	0.218	0.190	0.130	0.035	0.005	0.000	0.000	0.000	0.000	
	5	0.032	0.175	0.202	0.179	0.075	0.015	0.001	0.000	0.000	0.000	
	6	0.009	0.109	0.169	0.192	0.124	0.037	0.005	0.000	0.000	0.000	
	7	0.002	0.055	0.112	0.164	0.166	0.074	0.015	0.001	0.000	0.000	
	8	0.000	0.022	0.061	0.114	0.180	0.120	0.035	0.004	0.001	0.000	
	9	0.000	0.007	0.027	0.065	0.160	0.160	0.071	0.012	0.003	0.000	
	10	0.000	0.002	0.010	0.031	0.117	0.176	0.117	0.031	0.010	0.002	
	11	0.000	0.000	0.003	0.012	0.071	0.160	0.160	0.065	0.027	0.007	
	12	0.000	0.000	0.001	0.004	0.035	0.120	0.180	0.114	0.061	0.022	
	13	0.000	0.000	0.000	0.001	0.015	0.074	0.166	0.164	0.112	0.055	
	14	0.000	0.000	0.000	0.000	0.005	0.037	0.124	0.192	0.169	0.109	
	15	0.000	0.000	0.000	0.000	0.001	0.015	0.075	0.179	0.202	0.175	
	16	0.000	0.000	0.000	0.000	0.000	0.005	0.035	0.130	0.190	0.218	
	17	0.000	0.000	0.000	0.000	0.000	0.001	0.012	0.072	0.134	0.205	
	18	0.000	0.000	0.000	0.000	0.000	0.000	0.003	0.028	0.067	0.137	
	19	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.007	0.021	0.058	
	20	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.003	0.012	

Get More and Do More at Dummies.com®



Start with **FREE** Cheat Sheets

Cheat Sheets include

- Checklists
- Charts
- Common Instructions
- And Other Good Stuff!

To access the cheat sheet specifically for this book, go to
www.dummies.com/cheatsheet/statistics.

Get Smart at Dummies.com

Dummies.com makes your life easier with 1,000s of answers on everything from removing wallpaper to using the latest version of Windows.

Check out our

- Videos
- Illustrated Articles
- Step-by-Step Instructions

Plus, each month you can win valuable prizes by entering our Dummies.com sweepstakes.*

Want a weekly dose of Dummies? Sign up for Newsletters on

- Digital Photography
- Microsoft Windows & Office
- Personal Finance & Investing
- Health & Wellness
- Computing, iPods & Cell Phones
- eBay
- Internet
- Food, Home & Garden



*Sweepstakes not currently available in all countries; visit Dummies.com for official rules.

Find out "HOW" at Dummies.com

Making Everything Easier!™

Statistics II FOR DUMMIES®

Learn to:

- Increase your skills in data analysis
- Sort through and test models
- Make predictions
- Apply statistics to real-world situations

Deborah Rumsey, PhD

Author of *Statistics For Dummies* and
Statistics Workbook For Dummies



Statistics II For Dummies[®]

Visit www.dummies.com/cheatsheet/statistics2 to view this book's cheat sheet.

Table of Contents

[Introduction](#)

[About This Book](#)

[Conventions Used in This Book](#)

[What You're Not to Read](#)

[Foolish Assumptions](#)

[How This Book Is Organized](#)

[Part I: Tackling Data Analysis and Model-Building Basics](#)

[Part II: Using Different Types of Regression to Make Predictions](#)

[Part III: Analyzing Variance with ANOVA](#)

[Part IV: Building Strong Connections with Chi-Square Tests](#)

[Part V: Nonparametric Statistics: Rebels without a Distribution](#)

[Part VI: The Part of Tens](#)

[Icons Used in This Book](#)

[Where to Go from Here](#)

[Part I: Tackling Data Analysis and Model-Building Basics](#)

[Chapter 1: Beyond Number Crunching: The Art and Science of Data Analysis](#)

[Data Analysis: Looking before You Crunch](#)

[Nothing \(not even a straight line\) lasts forever](#)

[Data snooping isn't cool](#)

[No \(data\) fishing allowed](#)

[Getting the Big Picture: An Overview of Stats II](#)

[Population parameter](#)

[Sample statistic](#)

[Confidence interval](#)

Hypothesis test

Analysis of variance (ANOVA)

Multiple comparisons

Interaction effects

Correlation

Linear regression

Chi-square tests

Nonparametrics

Chapter 2: Finding the Right Analysis for the Job

Categorical versus Quantitative Variables

Statistics for Categorical Variables

Estimating a proportion

Comparing proportions

Looking for relationships between categorical variables

Building models to make predictions

Statistics for Quantitative Variables

Making estimates

Making comparisons

Exploring relationships

Predicting y using x

Avoiding Bias

Measuring Precision with Margin of Error

Knowing Your Limitations

Chapter 3: Reviewing Confidence Intervals and Hypothesis Tests

Estimating Parameters by Using Confidence Intervals

Getting the basics: The general form of a confidence interval

Finding the confidence interval for a population mean

What changes the margin of error?

Interpreting a confidence interval

What's the Hype about Hypothesis Tests?

What H_0 and H_a really represent

Gathering your evidence into a test statistic

Determining strength of evidence with a p-value

False alarms and missed opportunities: Type I and II errors

The power of a hypothesis test

Part II: Using Different Types of Regression to Make Predictions

Chapter 4: Getting in Line with Simple Linear Regression

Exploring Relationships with Scatterplots and Correlations

Using scatterplots to explore relationships

Collating the information by using the correlation coefficient

Building a Simple Linear Regression Model

Finding the best-fitting line to model your data

The y-intercept of the regression line

The slope of the regression line

Making point estimates by using the regression line

No Conclusion Left Behind: Tests and Confidence Intervals for Regression

Scrutinizing the slope

Inspecting the y-intercept

Building confidence intervals for the average response

Making the band with prediction intervals

Checking the Model's Fit (The Data, Not the Clothes!)

Defining the conditions

Finding and exploring the residuals

Using r^2 to measure model fit

Scoping for outliers

Knowing the Limitations of Your Regression Analysis

Avoiding slipping into cause-and-effect mode

Extrapolation: The ultimate no-no

Sometimes you need more than one variable

Chapter 5: Multiple Regression with Two X Variables

Getting to Know the Multiple Regression Model

Discovering the uses of multiple regression

Looking at the general form of the multiple regression model

Stepping through the analysis

Looking at x's and y's

Collecting the Data

Pinpointing Possible Relationships

Making scatterplots

Correlations: Examining the bond

Checking for Multicollinearity

Finding the Best-Fitting Model for Two x Variables

Getting the multiple regression coefficients

Interpreting the coefficients

Testing the coefficients

Predicting y by Using the x Variables

Checking the Fit of the Multiple Regression Model

Noting the conditions

Plotting a plan to check the conditions

Checking the three conditions

Chapter 6: How Can I Miss You If You Won't Leave? Regression Model Selection

Getting a Kick out of Estimating Punt Distance

Brainstorming variables and collecting data

Examining scatterplots and correlations

Just Like Buying Shoes: The Model Looks Nice, But Does It Fit?

Assessing the fit of multiple regression models

Model selection procedures

Chapter 7: Getting Ahead of the Learning Curve with Nonlinear Regression

Anticipating Nonlinear Regression

Starting Out with Scatterplots

Handling Curves in the Road with Polynomials

Bringing back polynomials

Searching for the best polynomial model

Using a second-degree polynomial to pass the quiz

Assessing the fit of a polynomial model

Making predictions

Going Up? Going Down? Go Exponential!

Recollecting exponential models

Searching for the best exponential model

Spreading secrets at an exponential rate

Chapter 8: Yes, No, Maybe So: Making Predictions by Using Logistic Regression

Understanding a Logistic Regression Model

How is logistic regression different from other regressions?

Using an S-curve to estimate probabilities

Interpreting the coefficients of the logistic regression model

The logistic regression model in action

Carrying Out a Logistic Regression Analysis

Running the analysis in Minitab

Finding the coefficients and making the model

Estimating p

Checking the fit of the model

Fitting the movie model

Part III: Analyzing Variance with ANOVA

Chapter 9: Testing Lots of Means? Come On Over to ANOVA!

Comparing Two Means with a t-Test

Evaluating More Means with ANOVA

Spitting seeds: A situation just waiting for ANOVA

Walking through the steps of ANOVA

Checking the Conditions

Verifying independence

Looking for what's normal

Taking note of spread

Setting Up the Hypotheses

Doing the F-Test

Running ANOVA in Minitab

Breaking down the variance into sums of squares

Locating those mean sums of squares

Figuring the F-statistic

Making conclusions from ANOVA

What's next?

Checking the Fit of the ANOVA Model

Chapter 10: Sorting Out the Means with Multiple Comparisons

Following Up after ANOVA

Comparing cellphone minutes: An example

Setting the stage for multiple comparison procedures

Pinpointing Differing Means with Fisher and Tukey

Fishing for differences with Fisher's LSD

Using Fisher's new and improved LSD

Separating the turkeys with Tukey's test

Examining the Output to Determine the Analysis

So Many Other Procedures, So Little Time!

Controlling for baloney with the Bonferroni adjustment

Comparing combinations by using Scheffe's method

Finding out whodunit with Dunnett's test

Staying cool with Student Newman-Keuls

Duncan's multiple range test

Going nonparametric with the Kruskal-Wallis test

Chapter 11: Finding Your Way through Two-Way ANOVA

Setting Up the Two-Way ANOVA Model

Determining the treatments

Stepping through the sums of squares

Understanding Interaction Effects

What is interaction, anyway?

Interacting with interaction plots

Testing the Terms in Two-Way ANOVA

Running the Two-Way ANOVA Table

Interpreting the results: Numbers and graphs

Are Whites Whiter in Hot Water? Two-Way ANOVA Investigates

Chapter 12: Regression and ANOVA: Surprise Relatives!

Seeing Regression through the Eyes of Variation

Spotting variability and finding an “x-planation”

Getting results with regression

Assessing the fit of the regression model

Regression and ANOVA: A Meeting of the Models

Comparing sums of squares

Dividing up the degrees of freedom

Bringing regression to the ANOVA table

Relating the F- and t-statistics: The final frontier

Part IV: Building Strong Connections with Chi-Square Tests

Chapter 13: Forming Associations with Two-Way Tables

Breaking Down a Two-Way Table

Organizing data into a two-way table

Filling in the cell counts

Making marginal totals

Breaking Down the Probabilities

Marginal probabilities

Joint probabilities

Conditional probabilities

Trying To Be Independent

Checking for independence between two categories

Checking for independence between two variables

Demystifying Simpson's Paradox

Experiencing Simpson's Paradox

Figuring out why Simpson's Paradox occurs

Keeping one eye open for Simpson's Paradox

Chapter 14: Being Independent Enough for the Chi-Square Test

The Chi-square Test for Independence

Collecting and organizing the data

Determining the hypotheses

Figuring expected cell counts

Checking the conditions for the test

Calculating the Chi-square test statistic

Finding your results on the Chi-square table

Drawing your conclusions

Putting the Chi-square to the test

Comparing Two Tests for Comparing Two Proportions

Getting reacquainted with the Z-test for two population proportions

Equating Chi-square tests and Z-tests for a two-by-two table

Chapter 15: Using Chi-Square Tests for Goodness-of-Fit (Your Data, Not Your Jeans)

Finding the Goodness-of-Fit Statistic

What's observed versus what's expected

Calculating the goodness-of-fit statistic

Interpreting the Goodness-of-Fit Statistic Using a Chi-Square

Checking the conditions before you start

The steps of the Chi-square goodness-of-fit test

Part V: Nonparametric Statistics: Rebels without a Distribution

Chapter 16: Going Nonparametric

Arguing for Nonparametric Statistics

No need to fret if conditions aren't met

The median's in the spotlight for a change

So, what's the catch?

Mastering the Basics of Nonparametric Statistics

Sign

Rank

Signed rank

Rank sum

Chapter 17: All Signs Point to the Sign Test and Signed Rank Test

Reading the Signs: The Sign Test

Testing the median

Estimating the median

Testing matched pairs

Going a Step Further with the Signed Rank Test

A limitation of the sign test

Stepping through the signed rank test

Losing weight with signed ranks

Chapter 18: Pulling Rank with the Rank Sum Test

Conducting the Rank Sum Test

Checking the conditions

Stepping through the test

Stepping up the sample size

Performing a Rank Sum Test: Which Real Estate Agent Sells Homes Faster?

Checking the conditions for this test

Testing the hypotheses

Chapter 19: Do the Kruskal-Wallis and Rank the Sums with the Wilcoxon

Doing the Kruskal-Wallis Test to Compare More than Two Populations

Checking the conditions

Setting up the test

Conducting the test step by step

Pinpointing the Differences: The Wilcoxon Rank Sum Test

Pairing off with pairwise comparisons

Carrying out comparison tests to see who's different

Examining the medians to see how they're different

Chapter 20: Pointing Out Correlations with Spearman's Rank

Pickin' On Pearson and His Precious Conditions

Scoring with Spearman's Rank Correlation

Figuring Spearman's rank correlation

Watching Spearman at work: Relating aptitude to performance

Part VI: The Part of Tens

Chapter 21: Ten Common Errors in Statistical Conclusions

Claiming These Statistics Prove . . .

It's Not Technically Statistically Significant, But . . .

Concluding That x Causes y

Assuming the Data Was Normal

Only Reporting "Important" Results

Assuming a Bigger Sample Is Always Better

It's Not Technically Random, But . . .

Assuming That 1,000 Responses Is 1,000 Responses

Of Course the Results Apply to the General Population

Deciding Just to Leave It Out

Chapter 22: Ten Ways to Get Ahead by Knowing Statistics

Asking the Right Questions

Being Skeptical

Collecting and Analyzing Data Correctly

Calling for Help

Retracing Someone Else's Steps

Putting the Pieces Together

Checking Your Answers

Explaining the Output

Making Convincing Recommendations

Establishing Yourself as the Statistics Go-To Guy or Gal

Chapter 23: Ten Cool Jobs That Use Statistics

Pollster

Ornithologist (Bird Watcher)

Sportscaster or Sportswriter

Journalist

Crime Fighter

Medical Professional

Marketing Executive

Lawyer

Stock Broker

Appendix: Reference Tables

Cheat Sheet

Statistics II FOR DUMMIES®

by Deborah Rumsey, PhD



Wiley Publishing, Inc.

Statistics II For Dummies®

Published by
John Wiley & Sons, Inc.
111 River St.
Hoboken, NJ 07030-5774

www.wiley.com

Copyright © 2009 by Wiley Publishing, Inc., Indianapolis, Indiana

Published by Wiley Publishing, Inc., Indianapolis, Indiana

Published simultaneously in Canada

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permissions>.

Trademarks: Wiley, the Wiley Publishing logo, For Dummies, the Dummies Man logo, A Reference for the Rest of Us!, The Dummies Way, Dummies Daily, The Fun and Easy Way, Dummies.com, Making Everything Easier, and related trade dress are trademarks or registered trademarks of John Wiley & Sons, Inc. and/or its affiliates in the United States and other countries, and may not be used without written permission. All other trademarks are the property of their respective owners. Wiley Publishing, Inc., is not associated with any product or vendor mentioned in this book.

Limit of Liability/Disclaimer of Warranty: The publisher and the author make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation warranties of fitness for a particular purpose. No warranty may be created or extended by sales or promotional materials. The advice and strategies contained herein may not be suitable for every situation. This work is sold with the understanding that the publisher is not engaged in rendering legal, accounting, or other professional services. If professional assistance is required, the services of a competent professional person should be sought. Neither the publisher nor the author shall be liable for damages arising herefrom. The fact that an organization or Website is referred to in this work as a citation and/or a potential source of further information does not mean that the author or the publisher endorses the information the organization or Website may provide or recommendations it may make. Further, readers should be aware that Internet Websites listed in this work may have changed or disappeared between when this work was written and when it is read.

For general information on our other products and services, please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993, or fax 317-572-4002.

For technical support, please visit www.wiley.com/techsupport.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Library of Congress Control Number: 2009928737

ISBN: 978-0-470-46646-9

Manufactured in the United States of America

10 9 8 7 6 5 4 3 2 1



Dedication

To my husband Eric: My sun rises and sets with you. To my son Clint: I love you up to the moon and back.

About the Author

Deborah Rumsey has a PhD in Statistics from The Ohio State University (1993), where she's a Statistics Education Specialist/Auxiliary Faculty Member for the Department of Statistics. Dr. Rumsey has been given the distinction of being named a Fellow of the American Statistical Association. She has also won the Presidential Teaching Award from Kansas State University. She's the author of *Statistics For Dummies*, *Statistics Workbook For Dummies*, and *Probability For Dummies* and has published numerous papers and given many professional presentations on the subject of statistics education. Her passions include being with her family, bird watching, getting more seat time on her Kubota tractor, and cheering the Ohio State Buckeyes on to another National Championship.

Author's Acknowledgments

Thanks again to Lindsay Lefevere and Kathy Cox for giving me the opportunity to write this book; to Natalie Harris and Chrissy Guthrie for their unwavering support and perfect chiseling and molding of my words and ideas; to Kim Gilbert, University of Georgia, for a thorough technical view; and to Elizabeth Rea and Sarah Westfall for great copy-editing. Special thanks to Elizabeth Stasny for guidance and support from day one; and to Joan Garfield for constant inspiration and encouragement.

Publisher's Acknowledgments

We're proud of this book; please send us your comments through our Dummies online registration form located at <http://dummies.custhelp.com>. For other comments, please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993, or fax 317-572-4002.

Some of the people who helped bring this book to market include the following:

Acquisitions, Editorial, and Media Development

Project Editors: Natalie Faye Harris, Chrissy Guthrie

Acquisitions Editors: Lindsay Lefevere, Kathy Cox

Copy Editors: Elizabeth Rea, Sarah Westfall

Assistant Editor: Erin Calligan Mooney

Editorial Program Coordinator: Joe Niesen

Technical Editor: Kim Gilbert

Editorial Manager: Christine Meloy Beck

Editorial Assistants: Jennette ElNaggar, David Lutton

Cover Photos: iStock

Cartoons: Rich Tennant (www.the5thwave.com)

Composition Services

Project Coordinator: Lynsey Stanford

Layout and Graphics: Carl Byers, Carrie Cesavice, Julie Trippetti, Christin Swinford, Christine Williams

Proofreaders: Melissa D. Buddendeck, Caitie Copple

Indexer: Potomac Indexing, LLC

Publishing and Editorial for Consumer Dummies

Diane Graves Steele, Vice President and Publisher, Consumer Dummies

Kristin Ferguson-Wagstaffe, Product Development Director, Consumer Dummies

Ensley Eikenburg, Associate Publisher, Travel

Kelly Regan, Editorial Director, Travel

Publishing for Technology Dummies

Andy Cummings, Vice President and Publisher, Dummies Technology/General User

Composition Services

Debbie Stailey, Director of Composition Services

Introduction

So you've gone through some of the basics of statistics. Means, medians, and standard deviations all ring a bell. You know about surveys and experiments and the basic ideas of correlation and simple regression. You've studied probability, margin of error, and a few hypothesis tests and confidence intervals. Are you ready to load your statistical toolbox with a new level of tools? *Statistics II For Dummies* picks up right where *Statistics For Dummies* (Wiley) leaves off and keeps you moving along the road of statistical ideas and techniques in a positive, step-by-step way.

The focus of *Statistics II For Dummies* is on finding more ways of analyzing data. I provide step-by-step instructions for using techniques such as multiple regression, nonlinear regression, one-way and two-way analysis of variance (ANOVA), Chi-square tests, and nonparametric statistics. Using these new techniques, you estimate, investigate, correlate, and congregate even more variables based on the information at hand.

About This Book

This book is designed for those who have completed the basic concepts of statistics through confidence intervals and hypothesis testing (found in *Statistics For Dummies*) and are ready to plow ahead to get through the final part of Stats I, or to tackle Stats II. However, I do pepper in some brief overviews of Stats I as needed, just to remind you of what was covered and make sure you're up to speed. For each new technique, you get an overview of when and why it's used, how to know when you need it, step-by-step directions on how to do it, and tips and tricks from a seasoned data analyst (yours truly). Because it's very important to be able to know which method to use when, I emphasize what makes each technique distinct and what the results say. You also see many applications of the techniques used in real life.

I also include interpretation of computer output for data analysis purposes. I show you how to use the software to get the results, but I focus more on how to interpret the results found in the output, because you're more likely to be interpreting this kind of information rather than doing the programming specifically. And because the equations and calculations can get too involved by hand, you often use a computer to get your results. I include instructions for using Minitab to conduct many of the calculations in this book. Most statistics teachers who cover these topics hold this philosophy as well. (What a relief!)

This book is different from the other Stats II books in many ways. Notably, this book features

- ✓ **Full explanations of Stats II concepts.** Many statistics textbooks squeeze all the Stats II topics at the very end of Stats I coverage; as a result, these topics tend to get condensed and presented as if they're optional. But no worries; I take the time to

clearly and fully explain all the information you need to survive and thrive.

- ✓ **Dissection of computer output.** Throughout the book, I present many examples that use statistical software to analyze the data. In each case, I present the computer output and explain how I got it and what it means.
- ✓ **An extensive number of examples.** I include plenty of examples to cover the many different types of problems you'll face.
- ✓ **Lots of tips, strategies, and warnings.** I share with you some trade secrets, based on my experience teaching and supporting students and grading their papers.
- ✓ **Understandable language.** I try to keep things conversational to help you understand, remember, and put into practice statistical definitions, techniques, and processes.
- ✓ **Clear and concise step-by-step procedures.** In most chapters, you can find steps that intuitively explain how to work through Stats II problems — and remember how to do it on your own later on.

Conventions Used in This Book

Throughout this book, I've used several conventions that I want you to be aware of:

- ✓ I indicate multiplication by using a times sign, indicated by a lowered asterisk, *.
- ✓ I indicate the null and alternative hypotheses as H_0 (for the null hypothesis) and H_a (for the alternative hypothesis).
- ✓ The statistical software package I use and display throughout the book is Minitab 14, but I simply refer to it as Minitab.
- ✓ Whenever I introduce a new term, I italicize it.
- ✓ Keywords and numbered steps appear in **boldface**.
- ✓ Web sites and e-mail addresses appear in monofont.

What You're Not to Read

At times I get into some of the more technical details of formulas and procedures for those individuals who may need to know about them — or just really want to get the full story. These minutiae are marked with a Technical Stuff icon. I also include sidebars as an aside to the essential text, usually in the form of a real-life statistics example or some bonus info you may find interesting. You can feel free to skip those icons and sidebars because you won't miss any of the main information you need (but by reading them, you may just be able to impress your stats professor with your above-and-beyond knowledge of Stats III!).

Foolish Assumptions

Because this book deals with Stats II, I assume you have one previous course in introductory statistics under your belt (or at least have read *Statistics For Dummies*), with topics taking you up through the Central Limit Theorem and perhaps an introduction to confidence intervals and hypothesis tests (although I review these concepts briefly in Chapter 3). Prior experience with simple linear regression isn't necessary. Only college algebra is needed for the mathematics details. And, some experience using statistical software is a plus but not required.

As a student, you may be covering these topics in one of two ways: either at the tail end of your Stats I course (perhaps in a hurried way, but in some way nonetheless); or through a two-course sequence in statistics in which the topics in this book are the focus of the second course. If so, this book provides you the information you need to do well in those courses.

You may simply be interested in Stats II from an everyday point of view, or perhaps you want to add to your understanding of studies and statistical results presented in the media. If this sounds like you, you can find plenty of real-world examples and applications of these statistical techniques in action as well as cautions for interpreting them.

How This Book Is Organized

This book is organized into five major parts that explore the main topic areas in Stats II, along with one bonus part that offers a series of quick top-ten references for you to use. Each part contains chapters that break down the part's major objective into understandable pieces. The nonlinear setup of this book allows you to skip around and still have easy access to and understanding of any given topic.

Part I: Tackling Data Analysis and Model-Building Basics

This part goes over the big ideas of descriptive and inferential statistics and simple linear regression in the context of model-building and decision-making. Some material from Stats I receives a quick review. I also present you with the typical jargon of Stats II.

Part II: Using Different Types of Regression to Make Predictions

In this part, you can review and extend the ideas of simple linear regression to the process of using more than one predictor variable. This part presents techniques for

dealing with data that follows a curve (nonlinear models) and models for yes or no data used to make predictions about whether or not an event will happen (logistic regression). It includes all you need to know about conditions, diagnostics, model-building, data-analysis techniques, and interpreting results.

Part III: Analyzing Variance with ANOVA

You may want to compare the means of more than two populations, and that requires that you use analysis of variance (ANOVA). This part discusses the basic conditions required, the *F*-test, one-way and two-way ANOVA, and multiple comparisons. The final goal of these analyses is to show whether the means of the given populations are different and if so, which ones are higher or lower than the rest.

Part IV: Building Strong Connections with Chi-Square Tests

This part deals with the Chi-square distribution and how you can use it to model and test categorical (qualitative) data. You find out how to test for independence of two categorical variables using a Chi-square test. (No more making speculations just by looking at the data in a two-way table!) You also see how to use a Chi-square to test how well a model for categorical data fits.

Part V: Nonparametric Statistics: Rebels without a Distribution

This part helps you with techniques used in situations where you can't (or don't want to) assume your data comes from a population with a certain distribution, such as when your population isn't normal (the condition required by most other methods in Stats II).

Part VI: The Part of Tens

Reading this part can give you an edge in a major area beyond the formulas and techniques of Stats II: ending the problem right (knowing what kinds of conclusions you can and can't make). You also get to know Stats II in the real world, namely how it can help you stand out in a crowd.

You also can find an appendix at the back of this book that contains all the tables you need to understand and complete the calculations in this book.

Icons Used in This Book

I use icons in this book to draw your attention to certain text features that occur on a regular basis. Think of the icons as road signs that you encounter on a trip. Some signs

tell you about shortcuts, and others offer more information that you may need; some signs alert you to possible warnings, while others leave you with something to remember.



When you see this icon, it means I'm explaining how to carry out that particular data analysis using Minitab. I also explain the information you get in the computer output so you can interpret your results.



I use this icon to reinforce certain ideas that are critical for success in Stats II, such as things I think are important to review as you prepare for an exam.



When you see this icon, you can skip over the information if you don't want to get into the nitty-gritty details. They exist mainly for people who have a special interest or obligation to know more about the more technical aspects of certain statistical issues.



This icon points to helpful hints, ideas, or shortcuts that you can use to save time; it also includes alternative ways to think about a particular concept.



I use warning icons to help you stay away from common misconceptions and pitfalls you may face when dealing with ideas and techniques related to Stats II.

Where to Go from Here

This book is written in a nonlinear way, so you can start anywhere and still understand what's happening. However, I can make some recommendations if you want some direction on where to start.

If you're thoroughly familiar with the ideas of hypothesis testing and simple linear regression, start with Chapter 5 (multiple regression). Use Chapter 1 if you need a reference for the jargon that statisticians use in Stats II.

If you've covered all topics up through the various types of regression (simple, multiple, nonlinear, and logistic) or a subset of those as your professor deemed important, proceed to Chapter 9, the basics of analysis of variance (ANOVA).

Chapter 14 is the place to begin if you want to tackle categorical (qualitative) variables before hitting the quantitative stuff. You can work with the Chi-square test there.

Nonparametric statistics are presented starting with Chapter 16. This area is a hot topic in today's statistics courses, yet it's also one that doesn't seem to get as much space in textbooks as it should. Start here if you want the full details on the most common nonparametric procedures.

Part I

Tackling Data Analysis and Model-Building Basics

The 5th Wave By Rich Tennant



In this part...

To get you up and moving from the foundational concepts of statistics (covered in your Stats I textbook as well as *Statistics For Dummies*) to the new and exciting methods presented in this book, I first go over the basics of data analysis, important terminology, main goals and concepts of model-building, and tips for choosing appropriate statistics to fit the job. I refresh your memory of the most heavily referred to items from Stats I, and you also get a head start on making and looking at some basic computer output.

Chapter 1

Beyond Number Crunching: The Art and Science of Data Analysis

In This Chapter

- ▶ Realizing your role as a data analyst
 - ▶ Avoiding statistical faux pas
 - ▶ Delving into the jargon of Stats II
-

Because you’re reading this book, you’re likely familiar with the basics of statistics and you’re ready to take it up a notch. That next level involves using what you know, picking up a few more tools and techniques, and finally putting it all to use to help you answer more realistic questions by using real data. In statistical terms, you’re ready to enter the world of the *data analyst*.

In this chapter, you review the terms involved in statistics as they pertain to data analysis at the Stats II level. You get a glimpse of the impact that your results can have by seeing what these analysis techniques can do. You also gain insight into some of the common misuses of data analysis and their effects.

Data Analysis: Looking before You Crunch

It used to be that statisticians were the only ones who really analyzed data because the only computer programs available were very complicated to use, requiring a great deal of knowledge about statistics to set up and carry out analyses. The calculations were tedious and at times unpredictable, and they required a thorough understanding of the theories and methods behind the calculations to get correct and reliable answers.

Today, anyone who wants to analyze data can do it easily. Many user-friendly statistical software packages are made expressly for that purpose — Microsoft Excel, Minitab, SAS, and SPSS are just a few. Free online programs are available, too, such as Stat Crunch, to help you do just what it says — crunch your numbers and get an answer.

Each software package has its own pros and cons (and its own users and protesters). My software of choice and the one I reference throughout this book is Minitab, because it’s very easy to use, the results are precise, and the software’s loaded with all the data-analysis techniques used in Stats II. Although a site license for Minitab isn’t cheap, the student version is available for rent for only a few bucks a semester.



The most important idea when applying statistical techniques to analyze data is to know what's going on behind the number crunching so you (not the computer) are in control of the analysis. That's why knowledge of Stats II is so critical.



Many people don't realize that statistical software can't tell you when to use and not to use a certain statistical technique. You have to determine that on your own. As a result, people think they're doing their analyses correctly, but they can end up making all kinds of mistakes. In the following sections, I give examples of some situations in which innocent data analyses can go wrong and why it's important to spot and avoid these mistakes before you start crunching numbers.

Bottom line: Today's software packages are too good to be true if you don't have a clear and thorough understanding of the Stats II that's underneath them.

Remembering the old days

In the old days, in order to determine whether different methods gave different results, you had to write a computer program using code that you had to take a class to learn. You had to type in your data in a specific way that the computer program demanded, and you had to submit your program to the computer and wait for the results. This method was time consuming and a general all-around pain.

The good news is that statistical software packages have undergone an incredible evolution in the last 10 to 15 years, to the point where you can now enter your data quickly and easily in almost any format. Moreover, the choices for data analysis are well organized and listed in pull-down menus. The results come instantly and successfully, and you can cut and paste them into a word-processing document without blinking an eye.

Nothing (not even a straight line) lasts forever

Bill Prediction is a statistics student studying the effect of study time on exam score. Bill collects data on statistics students and uses his trusty software package to predict exam score using study time. His computer comes up with the equation $y = 10x + 30$, where y represents the test score you get if you study a certain number of hours (x). Notice that this model is the equation of a straight line with a y -intercept of 30 and a slope of 10.

So Bill predicts, using this model, that if you don't study at all, you'll get a 30 on the exam (plugging $x = 0$ into the equation and solving for y ; this point represents the y -intercept of the line). And he predicts, using this model, that if you study for 5 hours, you'll get an exam score of $y = (10 * 5) + 30 = 80$. So, the point $(5, 80)$ is also on this line.

But then Bill goes a little crazy and wonders what would happen if you studied for 40

hours (since it always seems that long when he's studying). The computer tells him that if he studies for 40 hours, his test score is predicted to be $(10 * 40) + 30 = 430$ points. Wow, that's a lot of points! Problem is, the exam only goes up to a total of 100 points. Bill wonders where his computer went wrong.

But Bill puts the blame in the wrong place. He needs to remember that there are limits on the values of x that make sense in this equation. For example, because x is the amount of study time, x can never be a number less than zero. If you plug a negative number in for x , say $x = -10$, you get $y = (10 * -10) + 30 = -70$, which makes no sense. However, the equation itself doesn't know that, nor does the computer that found it. The computer simply graphs the line you give it, assuming it'll go on forever in both the positive and negative directions.



After you get a statistical equation or model, you need to specify for what values the equation applies. Equations don't know when they work and when they don't; it's up to the data analyst to determine that. This idea is the same for applying the results of any data analysis that you do.

Data snooping isn't cool



Statisticians have come up with a saying that you may have heard: "Figures don't lie. Liars figure." Make sure that you find out about all the analyses that were performed on a data set, not just the ones reported as being statistically significant.

Suppose Bill Prediction (from the previous section) decides to try to predict scores on a biology exam based on study time, but this time his model doesn't fit. Not one to give in, Bill insists there must be some other factors that predict biology exam scores besides study time, and he sets out to find them.

Bill measures everything from soup to nuts. His set of 20 possible variables includes study time, GPA, previous experience in statistics, math grades in high school, and whether you chew gum during the exam. After his multitude of various correlation analyses, the variables that Bill found to be related to exam score were study time, math grades in high school, GPA, and gum chewing during the exam. It turns out that this particular model fits pretty well (by criteria I discuss in Chapter 5 on multiple linear regression models).

But here's the problem: By looking at all possible correlations between his 20 variables and exam score, Bill is actually doing 20 separate statistical analyses. Under typical conditions that I describe in Chapter 3, each statistical analysis has a 5 percent chance of being wrong just by chance. I bet you can guess which one of Bill's correlations likely came out wrong in this case. And hopefully he didn't rely on a stick of gum to boost his

grade in biology.



Looking at data until you find something in it is called *data snooping*. Data snooping results in giving the researcher his five minutes of fame but then leads him to lose all credibility because no one can repeat his results.

No (data) fishing allowed

Some folks just don't take no for an answer, and when it comes to analyzing data, that can lead to trouble.

Sue Gonnafindit is a determined researcher. She believes that her horse can count by stomping his foot. (For example, she says "2" and her horse stumps twice.) Sue collects data on her horse for four weeks, recording the percentage of time the horse gets the counting right. She runs the appropriate statistical analysis on her data and is shocked to find no significant difference between her horse's results and those you would get simply by guessing.

Determined to prove her results are real, Sue looks for other types of analyses that exist and plugs her data into anything and everything she can find (never mind that those analyses are inappropriate to use in her situation). Using the famous hunt-and-peck method, at some point she eventually stumbles upon a significant result. However, the result is bogus because she tried so many analyses that weren't appropriate and ignored the results of the appropriate analysis because it didn't tell her what she wanted to hear.

Funny thing, too. When Sue went on a late night TV program to show the world her incredible horse, someone in the audience noticed that whenever the horse got to the correct number of stomps, Sue would interrupt him and say "Good job!" and the horse quit stomping. He didn't know how to count; all he knew to do was to quit stomping when she said "Good job!"



Redoing analyses in different ways in order to try to get the results you want is called *data fishing*, and folks in the stats biz consider it to be a major no-no. (However, people unfortunately do it all too often to verify their strongly held beliefs.) By using the wrong data analysis for the sake of getting the results you desire, you mislead your audience into thinking that your hypothesis is actually correct when it may not be.

Getting the Big Picture: An Overview of Stats II

Stats II is an extension of Stats I (introductory statistics), so the jargon follows suit and

the techniques build on what you already know. In this section, you get an introduction to the terminology you use in Stats II along with a broad overview of the techniques that statisticians use to analyze data and find the story behind it. (If you’re still unsure about some of the terms from Stats I, you can consult your Stats I textbook or see my other book, *Statistics For Dummies* (Wiley), for a complete rundown.)

Population parameter



A *parameter* is a number that summarizes the *population*, which is the entire group you’re interested in investigating. Examples of parameters include the mean of a population, the median of a population, or the proportion of the population that falls into a certain category.

Suppose you want to determine the average length of a cellphone call among teenagers (ages 13–18). You’re not interested in making any comparisons; you just want to make a good guesstimate of the average time. So you want to estimate a population parameter (such as the mean or average). The population is all cellphone users between the ages of 13 and 18 years old. The parameter is the average length of a phone call this population makes.

Sample statistic

Typically you can’t determine population parameters exactly; you can only estimate them. But all is not lost; by taking a *sample* (a subset of individuals) from the population and studying it, you can come up with a good estimate of the population parameter. A *sample statistic* is a single number that summarizes that subset.

For example, in the cellphone scenario from the previous section, you select a sample of teenagers and measure the duration of their cellphone calls over a period of time (or look at their cellphone records if you can gain access legally). You take the average of the cellphone call duration. For example, the average duration of 100 cellphone calls may be 12.2 minutes — this average is a statistic. This particular statistic is called the *sample mean* because it’s the average value from your sample data.

Many different statistics are available to study different characteristics of a sample, such as the proportion, the median, and standard deviation.

Confidence interval

A *confidence interval* is a range of likely values for a population parameter. A confidence interval is based on a sample and the statistics that come from that sample. The main reason you want to provide a range of likely values rather than a single number is that sample results vary.

For example, suppose you want to estimate the percentage of people who eat chocolate. According to the Simmons Research Bureau, 78 percent of adults reported eating chocolate, and of those, 18 percent admitted eating sweets frequently. What's missing in these results? These numbers are only from a single sample of people, and those sample results are guaranteed to vary from sample to sample. You need some measure of how much you can expect those results to move if you were to repeat the study.

This expected variation in your statistic from sample to sample is measured by the *margin of error*, which reflects a certain number of standard deviations of your statistic you add and subtract to have a certain confidence in your results (see Chapter 3 for more on margin of error). If the chocolate-eater results were based on 1,000 people, the margin of error would be approximately 3 percent. This means the actual percentage of people who eat chocolate in the entire population is expected to be 78 percent, \pm 3 percent (that is, between 75 percent and 81 percent).

Hypothesis test

A *hypothesis test* is a statistical procedure that you use to test an existing claim about the population, using your data. The claim is noted by H_0 (the null hypothesis). If your data support the claim, you fail to reject H_0 . If your data don't support the claim, you reject H_0 and conclude an alternative hypothesis, H_a . The reason most people conduct a hypothesis test is not to merely show that their data support an existing claim, but rather to show that the existing claim is false, in favor of the alternative hypothesis.

The Pew Research Center studied the percentage of people who turn to ESPN for their sports news. Its statistics, based on a survey of about 1,000 people, found that in 2000, 23 percent of people said they go to ESPN; in 2004, only 20 percent reported going to ESPN. The question is this: Does this 3 percent reduction in viewers from 2000 to 2004 represent a significant trend that ESPN should worry about?

To test these differences formally, you can set up a hypothesis test. You set up your null hypothesis as the result you have to believe without your study, H_0 = No difference exists between 2000 and 2004 data for ESPN viewership. Your alternative hypothesis (H_a) is that a difference is there. To run a hypothesis test, you look at the difference between your statistic from your data and the claim that has been already made about the population (in H_0), and you measure how far apart they are in units of standard deviations.

With respect to the example, using the techniques from Chapter 3, the hypothesis test shows that 23 percent and 20 percent aren't far enough apart in terms of standard deviations to dispute the claim (H_0). You can't say the percentage of viewers of ESPN in the entire population changed from 2000 to 2004.



As with any statistical analysis, your conclusions can be wrong just by chance,

because your results are based on sample data, and sample results vary. In Chapter 3 I discuss the types of errors that can be made in conclusions from a hypothesis test.

Analysis of variance (ANOVA)

ANOVA is the acronym for *analysis of variance*. You use ANOVA in situations where you want to compare the means of more than two populations. For example, you want to compare the lifetimes of four brands of tires in number of miles. You take a random sample of 50 tires from each group, for a total of 200 tires, and set up an experiment to compare the lifetime of each tire, and record it. You have four means and four standard deviations now, one for each data set.

Then, to test for differences in average lifetime for the four brands of tires, you basically compare the variability between the four data sets to the variability within the entire data set, using a ratio. This ratio is called the *F-statistic*. If this ratio is large, the variability between the brands is more than the variability within the brands, giving evidence that not all the means are the same for the different tire brands. If the *F*-statistic is small, not enough difference exists between the treatment means compared to the general variability within the treatments themselves. In this case, you can't say that the means are different for the groups. (I give you the full scoop on ANOVA plus all the jargon, formulas, and computer output in Chapters 9 and 10.)

Multiple comparisons

Suppose you conduct ANOVA, and you find a difference in the average lifetimes of the four brands of tire (see the preceding section). Your next questions would probably be, “Which brands are different?” and “How different are they?” To answer these questions, use multiple-comparison procedures.

A *multiple-comparison procedure* is a statistical technique that compares means to each other and finds out which ones are different and which ones aren’t. With this information, you’re able to put the groups in order from those with the largest mean to those with the smallest mean, realizing that sometimes two or more groups were too close to tell and are placed together in a group.

Many different multiple-comparison procedures exist to compare individual means and come up with an ordering in the event that your *F*-statistic does find that some difference exists. Some of the multiple-comparison procedures include Tukey’s test, LSD, and pairwise *t*-tests. Some procedures are better than others, depending on the conditions and your goal as a data analyst. I discuss multiple-comparison procedures in detail in Chapter 11.



Never take that second step to compare the means of the groups if the ANOVA procedure doesn't find any significant results during the first step. Computer software will never stop you from doing a follow-up analysis, even if it's wrong to do so.

Interaction effects

An *interaction effect* in statistics operates the same way that it does in the world of medicine. Sometimes if you take two different medicines at the same time, the combined effect is much different than if you were to take the two individual medications separately.



Interaction effects can come up in statistical models that use two or more variables to explain or compare outcomes. In this case you can't automatically study the effect of each variable separately; you have to first examine whether or not an interaction effect is present.

For example, suppose medical researchers are studying a new drug for depression and want to know how this drug affects the change in blood pressure for a low dose versus a high dose. They also compare the effects for children versus adults. It could also be that dosage level affects the blood pressure of adults differently than the blood pressure of children. This type of model is called a *two-way ANOVA model*, with a possible interaction effect between the two factors (age group and dosage level). Chapter 11 covers this subject in depth.

Correlation

The term *correlation* is often misused. Statistically speaking, the correlation measures the strength and direction of the linear relationship between two *quantitative variables* (variables that represent counts or measurements only).



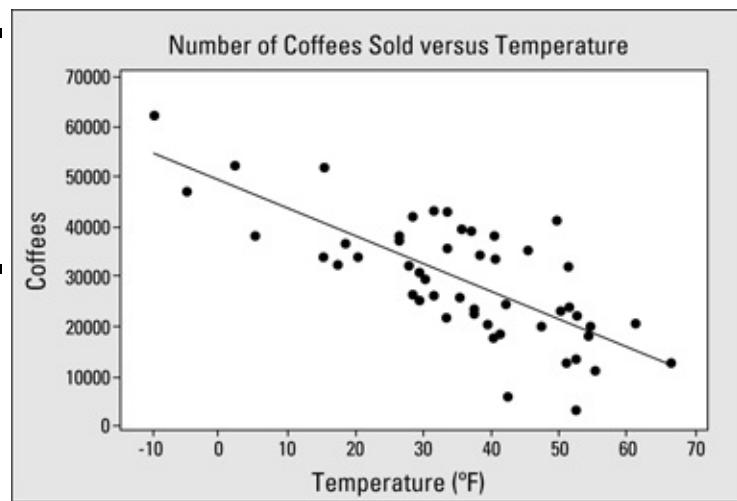
You aren't supposed to use correlation to talk about relationships unless the variables are quantitative. For example, it's wrong to say that a correlation exists between eye color and hair color. (In Chapter 14, you explore associations between two categorical variables.)

Correlation is a number between -1.0 and $+1.0$. A correlation of $+1$ indicates a perfect positive relationship; as you increase one variable, the other one increases in perfect sync. A correlation of -1.0 indicates a perfect negative relationship between the variables; as one variable increases, the other one decreases in perfect sync. A correlation of zero

means you found no linear relationship at all between the variables. Most correlations in the real world fall somewhere in between -1.0 and $+1.0$; the closer to -1.0 or $+1.0$, the stronger the relationship is; the closer to 0 , the weaker the relationship is.

Figure 1-1 shows a plot of the number of coffees sold at football games in Buffalo, New York, as well as the air temperature (in degrees Fahrenheit) at each game. This data set seems to follow a downhill straight line fairly well, indicating a negative correlation. The correlation turns out to be -0.741 ; number of coffees sold has a fairly strong negative relationship with the temperature of the football game. This makes sense because on days when the temperature is low, people get cold and want more coffee. I discuss correlation further, as it applies to model building, in Chapter 4.

Figure 1-1:
Coffees sold
at various air
temperatures
on football
game day.



Linear regression

After you've found a correlation and determined that two variables have a fairly strong linear relationship, you may want to try to make predictions for one variable based on the value of the other variable. For example, if you know that a fairly strong negative linear relationship exists between coffees sold and the air temperature at a football game (see the previous section), you may want to use this information to predict how much coffee is needed for a game, based on the temperature. This method of finding the best-fitting line is called *linear regression*.

Many different types of regression analyses exist, depending on your situation. When you use only one variable to predict the response, the method of regression is called *simple linear regression* (see Chapter 4). Simple linear regression is the best known of all the regression analyses and is a staple in the Stats I course sequence.

However, you use other flavors of regression for other situations.

- ✓ If you want to use more than one variable to predict a response, you use *multiple linear regression* (see Chapter 5).
- ✓ If you want to make predictions about a variable that has only two outcomes, yes or no, you use *logistic regression* (see Chapter 8).

- For relationships that don't follow a straight line, you have a technique called (no surprise) *nonlinear regression* (see Chapter 7).

Chi-square tests

Correlation and regression techniques all assume that the variable being studied in most detail (the response variable) is quantitative — that is, the variable measures or counts something. You can also run into situations where the data being studied isn't quantitative, but rather categorical — that is, the data represent categories, not measurements or counts. To study relationships in categorical data, you use a Chi-square test for independence. If the variables are found to be unrelated, they're declared independent. If they're found to be related, they're declared dependent.

Suppose you want to explore the relationship between gender and eating breakfast. Because each of these variables is categorical, or qualitative, you use a Chi-square test for independence. You survey 70 males and 70 females and find that 25 men eat breakfast and 45 do not; for the females, 35 do eat breakfast and 35 do not. Table 1-1 organizes this data and sets you up for the Chi-square test for this scenario.

Table 1-1 Table Setup for the Breakfast and Gender Question

	<i>Do Eat Breakfast</i>	<i>Don't Eat Breakfast</i>	<i>Total</i>
<i>Male</i>	25	45	70
<i>Female</i>	35	35	70



A Chi-square test first calculates what you expect to see in each cell of the table if the variables are independent (these values are brilliantly called the *expected cell counts*). The Chi-square test then compares these expected cell counts to what you observed in the data (called the *observed cell counts*) and compares them using a Chi-square statistic.

In the breakfast gender comparison, fewer males than females eat breakfast ($25 \div 70 = 35.7$ percent compared to $35 \div 70 = 50$ percent). Even though you know results will vary from sample to sample, this difference turns out to be enough to declare a relationship between gender and eating breakfast, according to the Chi-square test of independence. Chapter 14 reveals all the details of doing a Chi-square test.

You can also use the Chi-square test to see whether your theory about what percent of each group falls into a certain category is true or not. For example, can you guess what percentage of M&M'S fall into each color category? You can find more on these Chi-square variations, as well as the M&M'S question, in Chapter 15.

Nonparametrics

Nonparametrics is an entire area of statistics that provides analysis techniques to use when the conditions for the more traditional and commonly used methods aren't met. However, people sometimes forget or don't bother to check those conditions, and if the conditions are actually not met, the entire analysis goes out the window, and the conclusions go along with it!

Suppose you're trying to test a hypothesis about a population mean. The most common approach to use in this situation is a *t*-test. However, to use a *t*-test, the data needs to be collected from a population that has a normal distribution (that is, it has to have a bell-shaped curve). You collect data and graph it, and you find that it doesn't have a normal distribution; it has a skewed distribution. You're stuck — you can't use the common hypothesis test procedures you know and love (at least, you shouldn't use them).

This is where nonparametric procedures come in. Nonparametric procedures don't require nearly as many conditions be met as the regular parametric procedures do. In this situation of skewed data, it makes sense to run a hypothesis test for the median rather than the mean anyway, and plenty of nonparametric procedures exist for doing so.



If the conditions aren't met for a data-analysis procedure that you want to do, chances are that an equivalent nonparametric procedure is waiting in the wings. Most statistical software packages can do them just as easily as the regular (parametric) procedures.



Before doing a data analysis, statistical software packages don't automatically check conditions. It's up to you to check any and all appropriate conditions and, if they're seriously violated, to take another course of action. Many times a nonparametric procedure is just the ticket. For much more information on different nonparametric procedures, see Chapters 16 through 19.

Chapter 2

Finding the Right Analysis for the Job

In This Chapter

- ▶ Deciphering the difference between categorical and quantitative variables
 - ▶ Choosing appropriate statistical techniques for the task at hand
 - ▶ Evaluating bias and precision levels
 - ▶ Interpreting the results properly
-

One of the most critical elements of statistics and data analysis is the ability to choose the right statistical technique for each job. Carpenters and mechanics know the importance of having the right tool when they need it and the problems that can occur if they use the wrong tool. They also know that the right tool helps to increase their odds of getting the results they want the first time around, using the “work smarter, not harder” approach.

In this chapter, you look at some of the major statistical analysis techniques from the point of view of the carpenters and mechanics — knowing what each statistical tool is meant to do, how to use it, and when to use it. You also zoom in on mistakes some number crunchers make in applying the wrong analysis or doing too many analyses.



Knowing how to spot these problems can help you avoid making the same mistakes, but it also helps you to steer through the ocean of statistics that may await you in your job and in everyday life.

If many of the ideas you find in this chapter seem like a foreign language to you and you need more background information, don’t fret. Before continuing on in this chapter, head to your nearest Stats I book or check out another one of my books, *Statistics For Dummies* (Wiley).

Categorical versus Quantitative Variables

After you’ve collected all the data you need from your sample, you want to organize it, summarize it, and analyze it. Before plunging right into all the number crunching though, you need to first identify the type of data you’re dealing with. The type of data you have points you to the proper types of graphs, statistics, and analyses you’re able to use.



Before I begin, here's an important piece of jargon: Statisticians call any quantity or characteristic you measure on an individual a *variable*; the data collected on a variable is expected to vary from person to person (hence the creative name).

The two major types of variables are the following:

- ✓ **Categorical:** A *categorical variable*, also known as a *qualitative variable*, classifies the individual based on categories. For example, political affiliation may be classified into four categories: Democrat, Republican, Independent, and Other; gender as a variable takes on two possible categories: male and female. Categorical variables can take on numerical values only as placeholders.
- ✓ **Quantitative:** A *quantitative variable* measures or counts a quantifiable characteristic, such as height, weight, number of children you have, your GPA in college, or the number of hours of sleep you got last night. The quantitative variable value represents a quantity (count) or a measurement and has numerical meaning. That is, you can add, subtract, multiply, or divide the values of a quantitative variable, and the results make sense as numbers.

Because the two types of variables represent such different types of data, it makes sense that each type has its own set of statistics. Categorical variables, such as gender, are somewhat limited in terms of the statistics that can be performed on them.

For example, suppose you have a sample of 500 classmates classified by gender — 180 are male and 320 are female. How can you summarize this information? You already have the total number in each category (this statistic is called the *frequency*). You're off to a good start, but frequencies are hard to interpret because you find yourself trying to compare them to a total in your mind in order to get a proper comparison. For example, in this case you may be thinking, "One hundred and eighty males out of what? Let's see, it's out of 500. Hmm . . . what percentage is that?"

The next step is to find a means to relate these numbers to each other in an easy way. You can do this by using the *relative frequency*, which is the percentage of data that falls into a specific category of a categorical variable. You can find a category's relative frequency by dividing the frequency by the sample total and then multiplying by 100. In this case, you have $\frac{180}{500} = 0.36 * 100 = 36 \text{ percent males}$ and $\frac{320}{500} = 0.64 * 100 = 64 \text{ percent females}$.

You can also express the relative frequency as a proportion in each group by leaving the result in decimal form and not multiplying by 100. This statistic is called the *sample proportion*. In this example, the sample proportion of males is 0.36, and the sample proportion of females is 0.64.



You mainly summarize categorical variables by using two statistics — the number in each category (frequency) and the percentage (relative frequency) in each category.

Statistics for Categorical Variables

The types of statistics done on categorical data may seem limited; however, the wide variety of analyses you can perform using frequencies and relative frequencies offers answers to an extensive range of possible questions you may want to explore.

In this section, you see that the proportion in each group is the number-one statistic for summarizing categorical data. Beyond that, you see how you can use proportions to estimate, compare, and look for relationships between the groups that comprise the categorical data.

Estimating a proportion

You can use relative frequencies to make estimates about a single population proportion. (Refer to the earlier section “Categorical versus Quantitative Variables” for an explanation of relative frequencies.)

Suppose you want to know what proportion of females in the United States are Democrats. According to a sample of 29,839 female voters in the U.S. conducted by the Pew Research Foundation in 2003, the percentage of female Democrats was 36. Now, because the Pew researchers based these results on only a sample of the population and not on the entire population, their results will vary if they take another sample. This variation in sample results is cleverly called — you guessed it — sampling variability.

The sampling variability is measured by the *margin of error* (the amount that you add and subtract from your sample statistic), which for this sample is only about 0.5 percent. (To find out how to calculate margin of error, turn to Chapter 3.) That means that the estimated percentage of female Democrats in the U.S. voting population is somewhere between 35.5 percent and 36.5 percent.



The margin of error, combined with the sample proportion, forms what statisticians call a confidence interval for the population proportion. Recall from Stats I that a *confidence interval* is a range of likely values for a population parameter, formed by taking the sample statistic plus or minus the margin of error. (For more on confidence intervals, see Chapter 3.)

Comparing proportions

Researchers, the media, and even everyday folk like you and me love to compare groups (whether you like to admit it or not). For example, what proportion of Democrats support oil drilling in Alaska, compared to Republicans? What percentage of women watch college football, compared to men? What proportion of readers of *Statistics II For Dummies* pass their stats exams with flying colors, compared to nonreaders?

To answer these questions, you need to compare the sample proportions using a hypothesis test for two proportions (see Chapter 3 or your Stats I textbook).

Suppose you've collected data on a random sample of 1,000 voters in the U.S. and you want to compare the proportion of female voters to the proportion of male voters and find out whether they're equal. Suppose in your sample you find that the proportion of females is 0.53, and the proportion of males is 0.47. So for this sample of 1,000 people, you have a higher proportion of females than males.

But here's the big question: Are these sample proportions different enough to say that the entire population of American voters has more females in it than males? After all, sample results vary from sample to sample. The answer to this question requires comparing the sample proportions by using a hypothesis test for two proportions. I demonstrate and expand on this technique in Chapter 3.

Looking for relationships between categorical variables

Suppose you want to know whether two categorical variables are related; for example, is gender related to political affiliation? Answering this question requires putting the sample data into a two-way table (using rows and columns to represent the two variables) and analyzing the data by using a Chi-square test (see Chapter 14).

By following this process, you can determine if two categorical variables are independent (unrelated) or if a relationship exists between them. If you find a relationship, you can use percentages to describe it.

Table 2-1 shows an example of data organized in a two-way table. The data was collected by the Pew Research Foundation.

Table 2-1 Gender and Political Affiliation of 56,735 U.S. Voters

Gender	Republican	Democrat	Other
Males	32%	27%	41%
Females	29%	36%	35%

Notice that the percentage of male Republicans in the sample is 32 and the percentage of female Republicans in the sample is 29. These percentages are quite close in relative terms. However, the percentage of female Democrats seems much higher than the percentage of male Democrats (36 percent versus 27 percent); also, the percentage of males in the “Other” category is quite a bit higher than the percentage of females in the same category (41 percent versus 35 percent).

These large differences in the percentages indicate that gender and political affiliation are related in the sample. But do these trends carry over to the population of all American voters? This question requires a hypothesis test to answer. Because gender and political affiliation are both categorical variables, the particular hypothesis test you need in this situation is a Chi-square test. (I discuss Chi-square tests in detail in Chapter 14.)



To make a two-way table from a data set by using Minitab, first enter the data in two columns, where column one is the row variable (in this case, gender) and column two is the column variable (in this case, political affiliation). For example, suppose the first person is a male Democrat. In row one of Minitab, enter *M* (for male) in column one and *D* (Democrat) in column two. Then go to Stat>Tables>Cross Tabulation and Chi-square. Highlight column one and click Select to enter this variable in the For Rows line. Highlight column two and click Select to enter this variable in the For Columns line. Click OK.



People often use the word *correlation* to discuss relationships between variables, but in the world of statistics, correlation only relates to the relationship between two quantitative (numerical) variables, not two categorical variables. *Correlation* measures how closely the relationship between two quantitative variables, such as height and weight, follows a straight line and tells you the direction of that line as well. In total, for any two quantitative variables, *x* and *y*, the correlation measures the strength and direction of their linear relationship. As one increases, what does the other one do?

Because categorical variables don't have a numerical order to them, they don't increase or decrease in value. For example, just because male = 1 and female = 2 doesn't mean that a female is worth twice as much as a male (although some women may want to disagree). Therefore, you can't use the word *correlation* to describe the relationship between, say, gender and political affiliation. (Chapter 4 covers correlation.)

The appropriate term to describe the relationships of categorical variables is *association*. You can say that political affiliation is associated with gender and then explain how. (For full details on association, see Chapter 13.)

Building models to make predictions

You can build models to predict the value of a categorical variable based on other related information. In this case, building models is more than a lot of little plastic pieces and some irritatingly sticky glue.

When you build a statistical model, you look for variables that help explain, estimate, or predict some response you're interested in; the variables that do this are called

explanatory variables. You sort through the explanatory variables and figure out which ones do the best job of predicting the response. Then you put them together into a type of equation like $y = 2x + 4$ where x = shoe size and y = estimated calf length. That equation is a *model*.

For example, suppose you want to know which factors or variables can help you predict someone's political affiliation. Is a woman without children more likely to be a Republican or a Democrat? What about a middle-aged man who proclaims Hinduism as his religion?

In order for you to compare these complex relationships, you must build a model to evaluate each group's impact on political affiliation (or some other categorical variable). This kind of model-building is explored in-depth in Chapter 8, where I discuss the topic of logistic regression.



Logistic regression builds models to predict the outcome of a categorical variable, such as political affiliation. If you want to make predictions about a quantitative variable, such as income, you need to use the standard type of regression (check out Chapters 4 and 5).

Statistics for Quantitative Variables

Quantitative variables, unlike categorical variables, have a wider range of statistics that you can do, depending on what questions you want to ask. The main reason for this wider range is that *quantitative data* are numbers that represent measurements or counts, so it makes sense that you can order, add or subtract, and multiply or divide them — and the results all have numerical meaning. In this section, I present the major data-analysis techniques for quantitative data. I expand on each technique in later chapters of this book.

Making estimates

Quantitative variables take on numerical values that involve counts or measurements, so they have means, medians, standard deviations, and all those good things that categorical variables don't have. Researchers often want to know what the average or median value is for a population (these are called parameters). To do this requires taking a sample and making a good guess, also known as an estimate, of that parameter.

To find an estimate for any population parameter requires a confidence interval. For categorical variables, you would find a confidence interval to estimate the population mean, median, or standard deviation, but by far the most common parameter of interest is the population mean.

A confidence interval for the population mean is the sample mean plus or minus a margin of error. (To calculate the margin of error in this case, see Chapter 3.) The result will be a range of likely values you have produced for the real population mean. Because the variable is quantitative, the confidence interval will take on the same units as the variable does. For example, household incomes will be in thousands of dollars.

There is no rule of thumb regarding how large or small the margin of error should be for a quantitative variable; it depends on what the variable is counting or measuring. For example, if you want average household income for the state of New York, a margin of error of plus or minus \$5,000 is not unreasonable. If the variable is the average number of steps from the first floor to the second floor of a two-story home in the U.S., the margin of error will be much smaller. Estimates of categorical variables, on the other hand, are percentages; most people want those confidence intervals to be within plus or minus 2 to 3 percent.

Making comparisons

Suppose you want to look at income (a quantitative variable) and how it relates to a categorical variable, such as gender or region of the country. Your first question may be: Do males still make more money than females? In this case, you can compare the mean incomes of two populations — males and females. This assessment requires a hypothesis test of two means (often called a *t*-test for independent samples). I present more information on this technique in Chapter 3.



When comparing the means of *more* than two groups, don't simply look at all the possible *t*-tests that you can do on the pairs of means because you have to control for an overall error rate in your analysis. Too many analyses can result in errors — adding up to disaster. For example, if you conduct 100 hypothesis tests, each one with a 5 percent error rate, then 5 of those 100 tests will come out statistically significant on average, just by chance, even if no real relationship exists.

If you want to compare the average wage in different regions of the country (the East, the Midwest, the South, and the West, for example), this comparison requires a more sophisticated analysis because you're looking at four groups rather than just two. The procedure for comparing more than two means is called *analysis of variance* (ANOVA, for short), and I discuss this method in detail in Chapters 9 and 10.

Exploring relationships

One of the most common reasons data is collected is to look for relationships between variables. With quantitative variables, the most common type of relationship people look for is a linear relationship; that is, as one variable increases, does the other increase/decrease along with it in a similar way? Relationships between any variables are

examined using specialized plots and statistics. Since a linear relationship is so common, it has its own special statistic called correlation. You find out how statisticians make graphs and statistics to explore relationships in this section, paying particular attention to linear relationships.

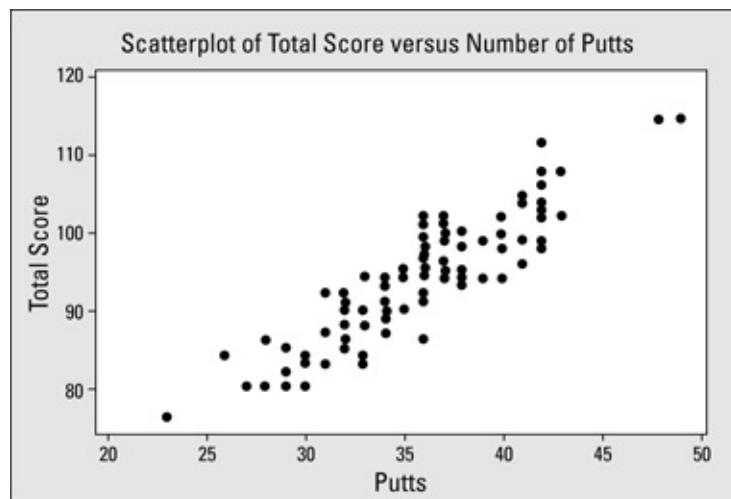
Suppose you're an avid golfer and you want to figure out how much time you should spend on your putting game. The question is this: Is the number of putts related to your total score? If the answer is yes, then spending time on your putting game makes sense. If not, then you can slack off on it a bit. Both of these variables are quantitative variables, and you're looking for a connection between them. You collect data on 100 rounds of golf played by golfers at your favorite course over a weekend. Following are the first few lines of your data set.

Round Number of Putts Total Score

1	23	76
2	27	80
3	28	80
4	29	80
5	30	80
6	29	82
7	30	83
8	31	83
9	33	83
10	26	84

The first step in looking for a connection between putts and total scores (or any other quantitative variables) is to make a scatterplot of the data. A *scatterplot* graphs your data set in two-dimensional space by using an X,Y plane. You can take a look at the scatterplot of the golf data in Figure 2-1. Here, *x* represents the number of putts, and *y* represents the total score. For example, the point in the lower-left corner of the graph represents someone who had only 23 putts and a total score of 75. (For instructions on making a scatterplot by using Minitab, see Chapter 4.)

Figure 2-1:
The two-dimensional scatterplot helps you look for relationships in data.



According to Figure 2-1, it appears that as the number of putts increases, so does the golfer's total score. It also shows that the variables increase in a linear way; that is, the data form a pattern that resembles a straight line. The relationship seems pretty strong — the number of putts plays a big part in determining the total score.

Now you need a measure of how strong the relationship is between x and y and whether it goes uphill or downhill. Different measures are used for different types of patterns seen in a scatterplot. Because the relationship we see in this case resembles a straight line, the correlation is the measure that we use to quantify the relationship. Correlation is the number that measures how close the points follow a straight line. Correlation is always between -1.0 and $+1.0$, and the more closely the points follow a straight line, the closer the correlation is to -1.0 or $+1.0$.

- ✓ **A positive correlation means that as x increases on the x -axis, y also increases on the y -axis.** Statisticians call this type of relationship an *uphill relationship*.
- ✓ **A negative correlation means that as x increases on the x -axis, y goes down.** Statisticians call this type of relationship — you guessed it — a *downhill relationship*.

For the golf data set, the correlation is $0.896 = 0.90$, which is extremely high as correlations go. The sign of the correlation is positive, so as you increase number of putts, your total score increases (an uphill relationship). For instructions on calculating a correlation in Minitab, see Chapter 4.

Predicting y using x

If you want to predict some response variable (y) using one explanatory variable (x) and you want to use a straight line to do it, you can use *simple linear regression* (see Chapter 4 for all the fine points on this topic). Linear regression finds the best-fitting line — called the *regression line* — that cuts through the data set. After you get the regression line, you can plug in a value of x and get your prediction for y . (For instructions on using Minitab to find the best-fitting line for your data, see Chapter 4.)

To use the golf example from the previous section, suppose you want to predict the total score you can get for a certain number of putts. In this case, you want to calculate the linear regression line. By running a regression analysis on the data set, the computer tells you that the best line to use to predict total score using number of putts is the following:

$$\text{Total score} = 39.6 + 1.52 * \text{Number of putts}$$

So if you have 35 putts in an 18-hole golf course, your total score is predicted to be about $39.6 + 1.52 * 35 = 92.8$, or 93. (Not bad for 18 holes!)



Don't try to predict y for x -values that fall outside the range of where the data was collected; you have no guarantee that the line still works outside of that range or that it will even make sense. For the golf example, you can't say that if x (the number of putts) = 1 the total score would be $39.6 + 1.52 * 1 = 41.12$ (unless you just call it good after your ball hits the green). This mistake is called *extrapolation*.

You can discover more about simple linear regression, and expansions on it, in Chapters 4 and 5.

Avoiding Bias

Bias is the bane of a statistician's existence; it's easy to create and very hard (if not impossible) to deal with in most situations. The statistical definition of *bias* is the systematic overestimation or underestimation of the actual value. In language the rest of us can understand, it means that the results are always off by a certain amount in a certain direction.

For example, a bathroom scale may always report a weight that's five pounds more than it should be (I'm convinced this is true of the scale at my doctor's office).

Bias can show up in a data set in a variety of different ways. Here are some of the most common ways bias can creep into your data:

- ✓ **Selecting the sample from the population:** Bias occurs when you either leave some groups out of the process that should have been included, or give certain groups too much weight.

For example, TV surveys that ask viewers to phone in their opinion are biased because no one has selected a prior sample of people to represent the population — viewers who want to be involved select themselves to participate by calling in on their own. Statisticians have found that folks who decide to participate in "call-in" or Web site polls are very likely to have stronger opinions than those who have been randomly selected but choose not to get involved in such polls. Such samples

are called *self-selected samples* and are typically very biased.

- ✓ **Designing the data-collection instrument:** Poorly designed instruments, including surveys and their questions, can result in inconsistent or even incorrect data. A survey question's wording plays a large role in whether or not results are biased. A leading question can make people feel like they should answer a certain way. For example, "Don't you think that the president should be allowed to have a line-item veto to prevent government spending waste?" Who would feel they should say *no* to that?
- ✓ **Collecting the data:** In this case, bias can infiltrate the results if someone makes errors in recording the data or if interviewers deviate from the script.
- ✓ **Deciding how and when the data is collected:** The time and place you collect data can affect whether your results are biased. For example, if you conduct a telephone survey during the middle of the day, people who work from 9 to 5 aren't able to participate. Depending on the issue, the timing of this survey could lead to biased results.

The best way to deal with bias is to avoid it in the first place, but you also can try to minimize it by

- ✓ **Using a random process to select the sample from the population.** The only way a sample is truly random is if every single member of the population has an equal chance of being selected. Self-selected samples aren't random.
- ✓ **Making sure the data is collected in a fair and consistent way.** Be sure to use neutral question wording and time the survey properly.

Don't put all your data in one basket!

An animal science researcher came to me one time with a data set he was so proud of. He was studying cows and the variables involved in helping determine their longevity. His super-mega data set contained over 100,000 observations. He was thinking, "Wow, this is gonna be great! I've been collecting this data for years and years, and I can finally have it analyzed. There's got to be loads of information I can get out of this. The papers I'll write, the talks I'll be invited to give . . . the raise I'll get!" He turned his precious data over to me with an expectant smile and sparkling eyes.

But after looking at his data for a few minutes I made a terrible realization — all his data came from exactly one cow. With no other cows to compare with and a sample size of just one, he had no way to even measure how much those results would vary if he wanted to apply them to another cow. His results were so biased toward that one animal that I couldn't do anything with the data. After I summoned the courage to tell him so, it took a while to peel him off the floor. The moral of the story, I suppose, is to run your big plans by a statistician before you go down a cow path like this guy did.

Measuring Precision with Margin of Error

Precision is the amount of movement you expect to have in your sample results if you repeat your entire study again with a new sample. Precision comes in two forms:

- ✓ **Low precision** means that you expect your sample results to move a lot (not a good thing).
- ✓ **High precision** means you expect your sample results to remain fairly close in the repeated samples (a good thing).

In this section, you find out what precision does and doesn't measure and you see how to measure the precision of a statistic in general terms.

Before you report or try to interpret any statistical results, you need to have some measurement of how much those results are expected to vary from sample to sample. This measurement is called the *margin of error*. You always hope, and may even assume, that statistical results shouldn't change much with another sample, but that's not always the case.

Up close and personal: Survey results

The Gallup Organization states its survey results in a universal, statistically correct format. Using a specific example from a recent survey it conducted, here's the language it uses to report its results:

"These results are based on telephone interviews with a randomly selected national sample of 1,002 adults, aged 18 years and older, conducted June 9–11, 2006. For results based on this sample, one can say with 95 percent confidence that the maximum error attributable to sampling and other random effects is ± 3 percentage points. In addition to sampling error, question wording and practical difficulties in conducting surveys can introduce error or bias into the findings of public opinion polls."

The first sentence of the quote refers to how the Gallup Organization collected the data, as well as the size of the sample. As you can guess, precision is related to the sample size, as seen in the section "Measuring Precision with Margin of Error."

The second sentence of the quote refers to the precision measurement: How much did Gallup expect these sample results to vary? The fact that Gallup is 95 percent confident means that if this process were repeated a large number of times, in 5 percent of the cases the results would be wrong, just by chance. This inconsistency occurs if the sample selected for the analysis doesn't represent the population — not due to biased reasons, but due to chance alone. Check out the section "Avoiding Bias" to get the info on why the third sentence is included in this quote.

The margin of error is affected by two elements:

- ✓ The sample size
- ✓ The amount of diversity in the population (also known as the population standard deviation)

You can read more about these elements in Chapter 3, but here's the big picture: As your sample size increases, you have more data to work with, and your results become more precise. As a result, the margin of error goes down.

On the other hand, a high amount of diversity in your population reduces your level of precision because the diversity makes it harder to get a handle on what's going on. As a result, the margin of error increases. (To offset this problem, just increase the sample size to get your precision back.)



To interpret the margin of error, just think of it as the amount of play you allow in your results to cover most of the other samples you could have taken.

Suppose you're trying to estimate the proportion of people in the population who support a certain issue, and you want to be 95 percent confident in your results. You sample 1,002 individuals and find that 65 percent support the issue. The margin of error for this survey turns out to be plus or minus 3 percentage points (you can find the details of this calculation in Chapter 3). That result means that you could expect the sample proportion of 65 percent to change by as much as 3 percentage points either way if you were to take a different sample of 1,002 individuals. In other words, you believe the actual population proportion is somewhere between $65 - 3 = 62$ percent and $65 + 3 = 68$ percent. That's the best you can say.



Any reported margin of error is calculated on the basis of having zero bias in the data. However, this assumption is rarely true. Before interpreting any margin of error, check first to be sure that the sampling process and the data-collection process don't contain any obvious sources of bias. Ignore results that are based on biased data, or at least take them with a great deal of skepticism.

For more details on how to calculate margin of error in various statistical techniques, turn to Chapter 3.

Knowing Your Limitations

The most important goal of any data analyst is to remain focused on the big picture — the question that you or someone else is asking — and make sure that the data analysis used is appropriate and comprehensive enough to answer that question correctly and fairly.



Here are some tips for analyzing data and interpreting the results, in terms of the

statistical procedures and techniques that you may use — at school, in your job, and in everyday life. These tips are implemented and reinforced throughout this book:

- ✓ **Be sure that the research question being asked is clear and definitive.** Some researchers don't want to be pinned down on any particular set of questions because they have the intent of mining the data — looking for any relationship they can find and then stating their results after the fact. This practice can lead to overanalyzing the data, making the results subject to skepticism by statisticians.
- ✓ **Double-check that you clearly understand the type of data being collected.** Is the data categorical or quantitative? The type of data used drives the approach that you take in the analysis.
- ✓ **Make sure that the statistical technique you use is designed to answer the research question.** If you want to make comparisons between two groups and your data is quantitative, use a hypothesis test for two means. If you want to compare five groups, use analysis of variance (ANOVA). Use this book as a resource to help you determine the technique you need.
- ✓ **Look for the limitations of the data analysis.** For example, if the researcher wants to know whether negative political ads affect the population of voters and she bases her study on a group of college students, you can find severe limitations here. For starters, student reactions to negative ads don't necessarily carry over to all voters in the population. In this case, it's best to limit the conclusions to college students in that class (which no researcher would ever want to do). Better to take a sample that represents the intended population of all voters in the first place (a much more difficult task, but well worth it).

Chapter 3

Reviewing Confidence Intervals and Hypothesis Tests

In This Chapter

- ▶ Utilizing confidence intervals to estimate parameters
 - ▶ Testing models by using hypothesis tests
 - ▶ Finding the probability of getting it right and getting it wrong
 - ▶ Discovering power in a large sample size
-

One of the major goals in statistics is to use the information you collect from a sample to get a better idea of what's going on in the entire population you're studying (because populations are generally large and exact info is often unknown). Unknown values that summarize the population are called *population parameters*. Researchers typically want to either get a handle on what those parameters are or test a hypothesis about the population parameters.

In Stats I, you probably went over confidence intervals and hypothesis tests for one and two population means and one and two population proportions. Your instructor hopefully emphasized that no matter which parameters you're trying to estimate or test, the general process is the same. If not, don't worry; this chapter drives that point home.

This chapter reviews the basic concepts of confidence intervals and hypothesis tests, including the probabilities of making errors by chance. I also discuss how statisticians measure the ability of a statistical procedure to do a good job — of detecting a real difference in the populations, for example.

Estimating Parameters by Using Confidence Intervals

Confidence intervals are a statistician's way of covering his you-know-what when it comes to estimating a population parameter. For example, instead of just giving a one-number guess as to what the average household income is in the United States, a statistician gives a range of likely values for this number. He does this because

- ✓ All good statisticians know sample results vary from sample to sample, so a one-number estimate isn't any good.
- ✓ Statisticians have developed some awfully nice formulas to give a range of likely values, so why not use them?

In this section, you get the general formula for a confidence interval, including the margin of error, and a good look at the common approach to building confidence intervals. I also discuss interpretation and the chance of making an error.

Getting the basics: The general form of a confidence interval

The big idea of a confidence interval is coming up with a range of likely values for a population parameter. The *confidence level* represents the chance that if you were to repeat your sample-taking over and over, you'd get a range of likely values that actually contains the actual population parameter. In other words, the confidence level is the long-term chance of being correct.



The general formula for a confidence interval is

$$\text{Confidence interval} = \text{Sample statistic} \pm \text{Margin of error}$$

The confidence interval has a certain level of precision (measured by the margin of error). Precision measures how close you expect your results to be to the truth.

For example, suppose you want to know the average amount of time a student at The Ohio State University spends listening to music on an MP3 player per day. The average time for the entire population of OSU students who are MP3-player users is the parameter you're looking for. You take a random sample of 1,000 students and find that the average time a student uses an MP3 player per day to listen to music is 2.5 hours, and the standard deviation is 0.5 hours. Is it right to say that the population of all OSU-student, MP3-player owners use their players an average of 2.5 hours per day for music listening? You hope and may assume that the average for the whole population is close to 2.5, but it probably isn't exact.

What's the solution to this problem? The solution is to not only report the average from your sample but along with it report some measure of how much you expect that sample average to vary from one sample to the next, with a certain level of confidence. The number that you use to represent this level of precision in your results is called the *margin of error*.

Finding the confidence interval for a population mean

The sample statistic part of the confidence-interval formula is fairly straightforward.

- ✓ **To estimate the population mean**, you use the sample mean plus or minus a margin of error, which is based on standard error. The sample mean has a standard error of $\frac{\sigma}{\sqrt{n}}$. In this formula, you can see the population standard deviation (σ) and

the sample size (n).

- ✓ To estimate the population proportion, you use the sample proportion plus or minus a margin of error.

In many cases, the standard deviation of the population, σ , is not known. To estimate the population mean by using a confidence interval when σ is unknown, you use the formula

$\bar{x} \pm t_{n-1} \left(\frac{s}{\sqrt{n}} \right)$. This formula contains the sample standard deviation (s), the sample size (n), and a t -value representing how many standard errors you want to add and subtract to get the confidence you need. To get the margin of error for the mean, you see the standard error, $\frac{s}{\sqrt{n}}$, is being multiplied by a factor of t . Notice that t has $n - 1$ as a subscript to indicate which of the myriad t -distributions you use for your confidence interval. The $n - 1$ is called *degrees of freedom*.

The value of t in this case represents the number of standard errors you add and subtract to or from the sample mean to get the confidence you want. If you want to be 95 percent confident, for example, you add and subtract 1.96 of those standard errors. If you want to be 99.7 percent confident, you add or subtract about three of them. (See Table A-1 in the appendix to find t -values for various confidence levels; use $\left(\frac{1 - \text{confidence level}}{2} \right)$ for the area to the right and find the t -value that goes with it.)



If you know the population standard deviation, you should certainly use it. In that case, you use the corresponding number from the Z-distribution (standard normal distribution) in the confidence interval formula. (The Z-distribution from your Stats I textbook can give you the numbers you need.) But I would be remiss in saying that while textbooks and teachers always include problems where σ is known, rarely is σ known in the real world. Why teach it this way? This issue is up for debate; for now just go with it, and I can keep you posted.

For the MP3 player example from the preceding section, a random sample of 1,000 all OSU students spends an average of 2.5 hours using their MP3 players to listen to music. The standard deviation is 0.5 hours. Plugging this information into the formula for a

confidence interval, you get $2.5 \pm 1.96 \left(\frac{0.5}{\sqrt{1,000}} \right)$. You conclude that all OSU-student MP3-owners spend an average of between 2.47 and 2.53 hours listening to music on their players.

What changes the margin of error?

What do you need to know in order to come up with a margin of error? Margin of error, in general, depends on three elements:

- ✓ The standard deviation of the population, σ (or an estimate of it, denoted by s , the sample standard deviation)

✓ The sample size, n

✓ The level of confidence you need

You can see these elements in action in the formula for margin of error of the sample mean: $\pm t_{n-1} * \frac{s}{\sqrt{n}}$. Here I assume that σ isn't known; t_{n-1} represents the value on the t -distribution table (see Table A-1 in the appendix) with $n - 1$ degrees of freedom.

Each of these three elements has a major role in determining how large the margin of error will be when you estimate the mean of a population. In the following sections, I show how each of the elements of the margin of error formula work separately and together to affect the size of the margin of error.

Population standard deviation

The standard deviation of the population is typically combined with the sample size in the margin of error formula, with the population standard deviation on top of the fraction and n on the bottom. (In this case, the standard error of the population, σ , is estimated by the standard deviation of the sample, s , because σ is typically unknown.)

This combination of standard deviation of the population and sample size is known as the *standard error* of your statistic. It measures how much the sample statistic deviates from its mean in the long term.



How does the standard deviation of the population (σ) affect margin of error? As it gets larger, the margin of error increases, so your range of likely values is wider.

Suppose you have two gas stations, one on a busy corner (gas station #1) and one farther off the main drag (gas station #2). You want to estimate the average time between customers at each station. At the busy gas station #1, customers are constantly using the gas pumps, so you basically have no downtime between customers. At gas station #2, customers sometimes come all at once, and sometimes you don't see a single person for an hour or more. So the time between customers varies quite a bit.

For which gas station would it be easier to estimate the overall average time between customers as a whole? Gas station #1 has much more consistency, which represents a smaller standard deviation of time between customers. Gas station #2 has much more variability in time between customers. That means σ for gas station #1 is smaller than σ for gas station #2. So the average time between customers is easier to estimate at gas station #1.

Sample size

Sample size affects margin of error in a very intuitive way. Suppose you're trying to

estimate the average number of pets per household in your city. Which sample size would give you better information: 10 homes or 100 homes? I hope you'd agree that 100 homes would give more precise information (as long as the data on those 100 homes was collected properly).

If you have more data to base your conclusions on and that data is collected properly, your results will be more precise. Precision is measured by margin of error, so as the sample size increases, the margin of error of your estimate goes down.



Bigger is only better in terms of sample size if the data is collected properly — that is, with minimal bias. If the quality of the data can't be maintained with a larger sample size, it does no good to have it.

Confidence level

For each problem at hand, you have to address how confident you need to be in your results over the long term, and, of course, more confidence comes with a price in the margin of error formula. This level of confidence in your results over the long term is reflected in a number called the *confidence level*, which you report as a percentage. In general, more confidence requires a wider range of likely values. So, as the confidence level increases, so does the margin of error.



Every margin of error is interpreted as plus or minus a certain number of standard errors. The number of standard errors added and subtracted is determined by the confidence level. If you need more confidence, you add and subtract more standard errors. If you need less confidence, you add and subtract fewer standard errors. The number that represents how many standard errors to add and subtract is different from situation to situation. For one population mean, you use a value on the *t*-distribution, represented by $tn - 1$, where n is the sample size (see Table A-1 in the appendix).

Suppose you have a sample size of 20, and you want to estimate the mean of a population with 90 percent confidence. The number of standard errors you add and subtract is represented by $tn - 1$, which in this case is $t_{19} = 1.73$. (To find these values of *t*, see Table A-1 in the appendix, with $n - 1$ degrees of freedom for the row, and $\frac{(1 - \text{confidence level})}{2}$ for the column.)

Now suppose you want to be 95 percent confident in your results, with the same sample size of $n = 20$. The degrees of freedom are $20 - 1 = 19$ (row) and the column is for $\frac{(1 - .95)}{2} = .025$. The t-table gives you the value of $t_{19} = 2.09$.

Notice that this value of t is larger than the value of t for 90 percent confidence, because in order to be more confident, you need to go out more standard deviations on the t -distribution table to cover more possible results.

Large confidence, narrow intervals — just the right size

A narrow confidence interval is much more desirable than a wide one. For example, claiming that the average cost of a new home is \$150,000 plus or minus \$100,000 isn't helpful at all because your estimate is anywhere between \$50,000 and \$250,000. (Who has an extra \$100,000 to throw around?) But you *do* want a high confidence level, so your statistician has to add and subtract more standard errors to get there, which makes the interval that much wider (a downer).

Wait, don't panic — you can have your cake and eat it too! If you know you want to have a high level of confidence but you don't want a wide confidence interval, just increase your sample size to meet that level of confidence.

Suppose the standard deviation of the house prices from a previous study is $s = \$15,000$, and you want to be 95 percent confident in your estimate of average house price. Using a large sample size, your value of t (from Table A-1 in the appendix) is 1.96.

With a sample of 100 homes, your margin of error is $\pm 1.96 \times \frac{15,000}{\sqrt{100}} = \$2,940$.

If this is too large for you but you still want 95 percent confidence, crank up your value of n . If you sample 500 homes, the margin of error decreases to $\pm 1.96 \times \frac{15,000}{\sqrt{500}}$, which brings you down to \$1,314.81.



You can use a formula to find the sample size you need to meet a desired margin

of error. That formula is $n = \left(\frac{t_{n-1} s}{MOE} \right)^2$, where MOE is the desired margin of error (as a proportion), s is the sample standard deviation, and t is the value on the t -distribution that corresponds with the confidence level you want. (For large sample sizes, the t -distribution is approximately equal to the Z -distribution; you can use the last line of Table A-1 in the appendix for the appropriate t -values, or use a Z -table from your Stats I textbook.)

Interpreting a confidence interval

Interpreting a confidence interval involves a couple of subtle but important issues. The big idea is that a *confidence interval* presents a range of likely values for the population parameter, based on your sample. However, you interpret it not in terms of your own sample, but in terms of an infinite number of other samples out there that could have been selected, yours just being one of them. For example, suppose 1,000 people each

took a sample and they each formed a 95 percent confidence interval for the mean. The “95 percent confidence” part means that of those 1,000 confidence intervals, about 950 of them can be expected to be correct on average. (Correct means the confidence interval actually contains the true value of the parameter.)



A 95 percent confidence interval doesn’t mean that your particular confidence interval has a 95 percent chance of capturing the actual value of the parameter; after the sample has been taken, the parameter is either in the interval or it isn’t. A confidence interval represents the chances of capturing the actual value of the population parameter over many different samples.

Suppose a polling organization wants to estimate the percentage of people in the United States who drive a car with more than 100,000 miles on it, and it wants to be 95 percent confident in its results. The organization takes a random sample of 1,200 people and finds that 420 of them (35 percent) drive a car with that minimum mileage; the margin of error turns out to be plus or minus 3 percent. (See your Stats I text for determining margin of error for percentages.)

The meaty part of the interpretation lies in the confidence level — in this case, the 95 percent. Because the organization took a sample of 1,200 people in the U.S., asked each of them whether his or her car has more than 100,000 miles on it, and made a confidence interval out of the results, the polling organization is, in essence, accounting for all the other samples out there that it could have gotten by building in the margin of error (± 3 percent). The organization wants to cover its bases on 95 percent of those other situations, and ± 3 percent satisfies that.

Another way of thinking about the confidence interval is to say that if the organization sampled 1,200 people over and over again and made a confidence interval from its results each time, 95 percent of those confidence intervals would be right. (You just have to hope that yours is one of those right results.)



Using stat notation, you can write confidence levels as $(1 - \alpha)\%$. So if you want 95 percent confidence, you write it as $1 - 0.05$. Here, α represents the chance that your confidence interval is one of the wrong ones. This number, α , is also related to the random chance of making a certain kind of error with a hypothesis test, which I explain in the later section “False alarms and missed opportunities: Type I and II errors.”

What’s the Hype about Hypothesis Tests?

Suppose a shipping company claims that its packages are on time 92 percent of the time,

or a campus official claims that 75 percent of students live off campus. If you're questioning these claims, how can you use statistics to investigate?

In this section, you see the big ideas of hypothesis testing that are the basis for the data-analysis techniques in this book. You review and expand on the concepts involved in a hypothesis test, including the hypotheses, the test statistic, and the p -value.

What H_0 and H_a really represent

You use a hypothesis test in situations where you have a certain model in mind and want to see whether that model fits your data. Your model may be one that just revolves around the population mean (testing whether that mean is equal to ten, for example). Your model may be testing the slope of a regression line (whether or not it's zero, for example, with zero meaning you find no relationship between x and y). You may be trying to use several different variables to predict the marketability of a product, and you believe a model using customer age, price, and shelf location can help predict it, so you need to run one or more hypothesis tests to see whether that model works. (This particular process is called multiple regression, and you can find more info on it in Chapter 5.)

A hypothesis test is made up of two hypotheses:

- ✓ **The null hypothesis, H_0 :** H_0 symbolizes the current situation — the one that everyone assumed was true until you got involved.
- ✓ **The alternative hypothesis, H_a :** H_a represents the alternative model that you want to consider. It stands for the researcher's hypothesis, and the burden of proof lies on the researcher.



H_0 is the model that's on trial. If you get enough evidence against it, you conclude H_a , which is the model you're claiming is the right one. If you don't get enough evidence against H_0 , then you can't say that your model (H_a) is the right one.

Gathering your evidence into a test statistic

A *test statistic* is the statistic from your sample, standardized so you can look it up on a table, basically. Although each hypothesis test is a little different, the main thought is the same. Take your statistic and standardize it in the appropriate way so you can use the corresponding table for it. Then look up your test statistic on a table to see where it stands. That table may be the t -table (Table A-1 in the appendix), the Chi-square table (Table A-3 in the appendix), or a different table. The type of test you need to use on your data dictates which table you use.

In the case of testing a hypothesis for a population mean, μ , you use the sample mean, \bar{x} ,

as your statistic. To standardize it, you take \bar{x} and convert it to a value of t by using the

$$t_{n-1} = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

formula where μ_0 is the value in H_0 .

This value is your test statistic, which you compare to the t -distribution.

Determining strength of evidence with a p-value

If you want to know whether your data has the brawn to stand up against H_0 , you need to figure out the p -value and compare it to a predetermined cutoff, α (typically 0.05). The p -value is a measure of the strength of your evidence against H_0 . You calculate the p -value through these steps:

1. Calculate the test statistic (refer to the preceding section for more info on this).
2. Look up the test statistic on the appropriate table (such as the t -table, Table A-1 in the appendix).
3. Find the percentage of values on the table that fall beyond your test statistic. This percentage is the p -value.
4. If your H_a is “not equal to,” double the percentage that you got in step three because your test statistic could have gone either way before the data was collected. (See your Stats I textbook or *Statistics For Dummies* for full details on obtaining p -values for hypothesis tests.)



Your friend α is the cutoff for your p -value. (α is typically set at 0.05, but sometimes it's 0.10.) If your p -value is less than your predetermined value of α , reject H_0 because you have sufficient evidence against it. If your p -value is greater than or equal to α , you can't reject H_0 .

For example, if your p -value is 0.002, your test statistic is so far away from H_0 that the chance of getting this result by chance is only 2 out of 1,000. So, you conclude that H_0 is very likely to be false. If your p -value turns out to be 0.30, this same result is expected to happen 30 percent of the time anyway, so you see no red flags there, and you can't reject H_0 . You don't have enough evidence against it. If your p -value is close to the cutoff line, say $p = 0.049$ or 0.51 , you say the result is marginal and let the reader make her own conclusions. That's the main advantage of the p -value: It lets other folks determine whether your evidence is strong enough to reject H_0 in their minds.

False alarms and missed opportunities: Type I and II errors

Any technique you use in statistics to make a conclusion about a population based on a sample of data has the chance of making an error. The errors I am talking about, Type I

and Type II errors, are due to random chance.



The way you set up your test can help to reduce these kinds of errors, but they're always out there. As a data analyst, you need to know how to measure and understand the impact of the errors that can occur with a hypothesis test and what you can do to possibly make those errors smaller. In the following sections, I show you how you can do just that.

Making false alarms with Type I errors

A *Type I error* is the conditional probability of rejecting H_0 , given that H_0 is true. I think of a Type I error as a false alarm: You blew the whistle when you shouldn't have.

The chance of making a Type I error is equal to α , which is predetermined before you begin collecting your data. This α is the same α that represents the chance of missing the boat in a confidence interval. It makes some sense that these two probabilities are both equal because the probability of rejecting H_0 when you shouldn't (a Type I error) is the same as the chance that the true population parameter falls out of the range of likely values when it shouldn't. That chance is α .

Suppose someone claims that the mean time to deliver packages for a company is 3.0 days on average (so H_0 is $\mu = 3.0$), but you believe it's not equal to that (so H_a is $\mu \neq 3.0$). Your α level is 0.05, and because you have a two-sided test, you have 0.025 on each side. Your sample of 100 packages has a mean of 3.5 days with a standard deviation of 1.5 days.

$$\frac{3.5 - 3.0}{1.5} = 3.33$$

The test statistic equals $\frac{3.5 - 3.0}{\sqrt{100}} = 0.5$, which is greater than 1.96 (the value on the last row and the 0.025 column of the *t*-distribution table — see Table A-1 in the appendix). So 3.0 is not a likely value for the mean time of delivery for all packages, and you reject H_0 .

But suppose that just by chance, your sample contained some longer than normal delivery times and that, in reality, the company's claim is right. You just made a Type I error. You made a false alarm about the company's claim.



To reduce the chance of a Type I error, reduce your value of α . However I don't recommend reducing it too far. On the positive side, this reduction makes it harder to reject H_0 because you need more evidence in your data to do so. On the negative side, by reducing your chance of a false alarm (Type I error) you increase the chance of a missed opportunity (a Type II error.)

Missing an opportunity with a Type II error

A *Type II error* is the conditional probability of not rejecting H_0 , given that H_0 is false. I

call it a missed opportunity because you were supposed to be able to find a problem with H_0 and reject it, but you didn't. You didn't blow the whistle when you should have.

The chance of making a Type II error depends on a couple of things:

- ✓ **Sample size:** If you have more data, you're less likely to miss something that's going on. For example, if a coin actually is unfair, flipping the coin only ten times may not reveal the problem. But if you flip the coin 1,000 times, you have a good chance of seeing a pattern that favors heads over tails, or vice versa.
- ✓ **Actual value of the parameter:** A Type II error is also related to how big the problem is that you're trying to uncover. For example, suppose a company claims that the average delivery time for packages is 3.5 days. If the actual average delivery time is 5.0 days, you won't have a very hard time detecting that with your sample (even a small sample). But if the actual average delivery time is 4.0 days, you have to do more work to actually detect the problem.



To reduce the chance of a Type II error, take a larger sample size. A greater sample size makes it easier to reject H_0 but increases the chance of a Type I error.



Type I and Type II errors sit on opposite ends of a seesaw — as one goes up, the other goes down. Try to meet in the middle by choosing a large sample size (the bigger, the better; see Figures 3-1 and 3-2) and a small α level (0.05 or less) for your hypothesis test.

The power of a hypothesis test

Type II errors, which I explain in the preceding section, show the downside of a hypothesis test. But statisticians, despite what many may think, actually try to look on the bright side once in a while; so instead of looking at the chance of *missing* a difference from H_0 that actually is there, they look at the chance of *detecting* a difference that really is there. This detection is called the *power of a hypothesis test*.



The power of a hypothesis test is $1 -$ the probability of making a Type II error. So *power* is a number between 0 and 1 that represents the chance that you rejected H_0 when H_0 was false. (You can even sing about it: "If H_0 is false and you know it, clap your hands. . . .") Remember that power (just like Type II errors) depends on two elements: the sample size and the actual value of the parameter (see the preceding section for a description of these elements).

In the following sections, you discover what power means in statistics (not being one of

the bigwigs, mind you); you also find out how to quantify power by using a power curve.

Throwing a power curve

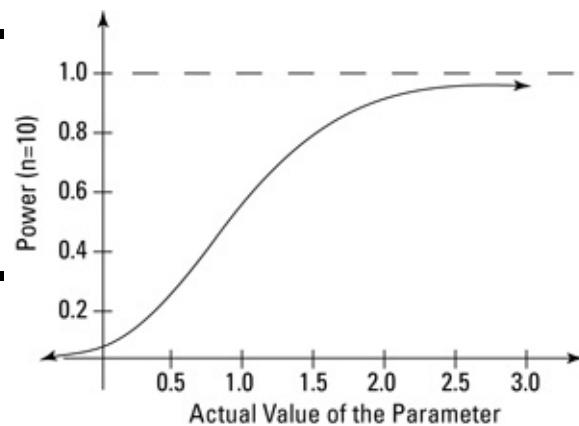
The specific calculations for the power of a hypothesis test are beyond the scope of this book (so you can take a sigh of relief), but computer programs and graphs are available online to show you what the power is for different hypothesis tests and various sample sizes (just type “power curve for the [blah blah blah] test” into an Internet search engine).

These graphs are called *power curves* for a hypothesis test. A power curve is a special kind of graph that gives you an idea of how much of a difference from H_0 you can detect with the sample size that you have. Because the precision of your test statistic increases as your sample size increases, sample size is directly related to power. But it also depends on how much of a difference from H_0 you’re trying to detect. For example, if a package delivery company claims that its packages arrive in 2 days or less, do you want to blow the whistle if it’s actually 2.1 days? Or wait until it’s 3 days? You need a much larger sample size to detect the 2.1-days situation versus the 3-days situation just because of the precision level needed.

In Figure 3-1, you can see the power curve for a particular test of $H_0: \mu = 0$ versus $H_a: \mu > 0$. You can assume that σ (the standard deviation of the population) is equal to 2 (I give you this value in each problem) and doesn’t change. I set the sample size at 10 throughout.

The horizontal (x) axis on the power curve shows a range of actual values of μ . For example, you hypothesize that μ is equal to 0, but it may actually be 0.5, 1.0, 2.0, 3.0, or any other possible value. If μ equals 0, then H_0 is true, and the chance of detecting this (and therefore rejecting H_0) is equal to 0.05, the set value of α . You work from that baseline. (Notice the low power in this situation makes sense because there’s nothing to detect for values of μ that are close to 0.) So, on the graph in Figure 3-1, when $x = 0$, you get a y -value of 0.05.

Figure 3-1:
Power curve
for $H_0: \mu = 0$
versus $H_a: \mu > 0$, for $n = 10$
and $\sigma = 2$.



Suppose that μ is actually 0.5, not 0, as you hypothesized. A computer tells you that the chance of rejecting H_0 (what you’re supposed to do here) is $0.197 - 0.20$, which is the power. So, you have about a 20-percent chance of detecting this difference with a sample

size of 10. As you move to the right, away from 0 on the horizontal (x) axis, you can see that the power goes up and the y -values get closer and closer to 1.0.

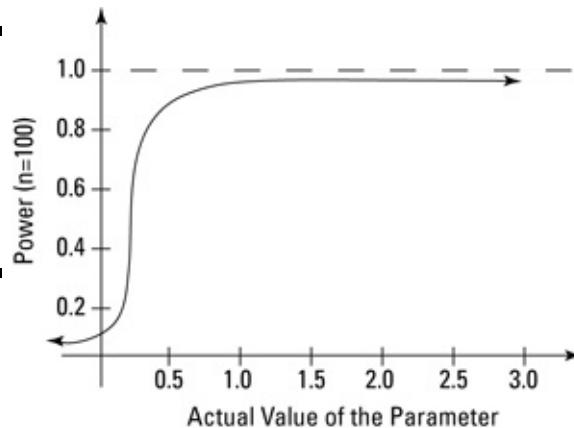
For example, if the actual value of μ is 1.0, the difference from 0 is easier to detect than if it's 0.50. In fact, the power at 1.0 is equal to $0.475 = 0.48$, so you have almost a 50 percent chance of catching the difference from H_0 in this case. And as the values of the mean increase, the power gets closer and closer to 1.0. Power never reaches 1.0 because statistics can never prove anything with 100 percent accuracy, but you can get close to 1.0 if the actual value is far enough from your hypothesis.

Controlling the sample size

How can you increase the power of your hypothesis test? You don't have any control over the actual value of the parameter, because that number is unknown. So what do you have control over? The sample size. As the sample size increases, it becomes easier to detect a real difference from H_0 .

Figure 3-2 shows the power curve with the same numbers as Figure 3-1, except for the sample size (n), which is 100 instead of 10. Notice that the curve increases much more quickly and approaches 1.0 when the actual mean is 1.0, compared to your hypothesis of 0. You want to see this kind of curve that moves up quickly toward the value of 1.0, while the actual values of the parameter increase on the x -axis.

Figure 3-2:
Power curve
for $H_0: \mu = 0$
versus $H_a: \mu > 0$, for $n = 100$ and $\sigma = 2$.



If you compare the power of your test when μ is 1.0 for the $n = 10$ situation (in Figure 3-1) versus the $n = 100$ situation (in Figure 3-2), you see that the power increases from 0.475 to more than 0.999. Table 3-1 shows the different values of power for the $n = 10$ case versus the $n = 100$ case, when you test $H_0: \mu = 0$ versus $H_a: \mu > 0$, assuming a value of $\sigma = 2$.

Table 3-1 Comparing the Values of Power for $n = 10$ Versus $n = 100$ (H_0 is $\mu = 0$)

Actual Value of μ	Power When $n = 10$	Power When $n = 100$
0.00	$0.050 = 0.05$	$0.050 = 0.05$
0.50	$0.197 = 0.20$	$0.804 = 0.81$
1.00	$0.475 = 0.48$	approx. 1.0
1.50	$0.766 = 0.77$	approx. 1.0

2.00

0.935 = 0.94

approx. 1.0

3.00

0.999 = approx. 1.0

approx. 1.0



You can find power curves for a variety of hypothesis tests under many different scenarios. Each has the same general look and feel to it: starting at the value of α when H_0 is true, increasing in an S-shape as you move from left to right on the x -axis, and finally approaching the value of 1.0 at some point. Power curves with large sample sizes approach 1.0 faster than power curves with low sample sizes.



It's possible to have too much power. For example, if you make the power curve for $n = 10,000$ and compare it to Figures 3-1 and 3-2, you find that it's practically at 1.0 already for any number other than 0.0 for the mean. In other words, the actual mean could be 0.05 and with your hypothesis $H_0: \mu = 0.00$, you would reject H_0 because of your huge sample size. Unless a researcher really wants to detect very small differences from H_0 (such as in medical studies or quality control situations), inflated values of n are usually suspect. People sometimes increase n just to be able to say they've found a difference, no matter how small, so watch for that. If you zoom in enough, you can always detect something, even if that something makes no practical difference. Beware of surveys and experiments with an excessive sample size, such as one in the tens of thousands. Their results are guaranteed to be inflated.

Power in manufacturing and medicine

The power of a test plays a role in the manufacturing process. Manufacturers often have very strict specifications regarding the size, weight, and/or quality of their products. During the manufacturing process, manufacturers want to be able to detect deviations from these specifications, even small ones, so they must determine how much of a difference from H_0 they want to detect, and then figure out the sample size needed in order to detect that difference when it appears. For example, if the candy bar is supposed to weigh 2.0 ounces, the manufacturer may want to blow the whistle if the actual average weight shifts to 2.2 ounces. Statisticians can work backward in calculating the power and find the sample size they need to know to stop the process.

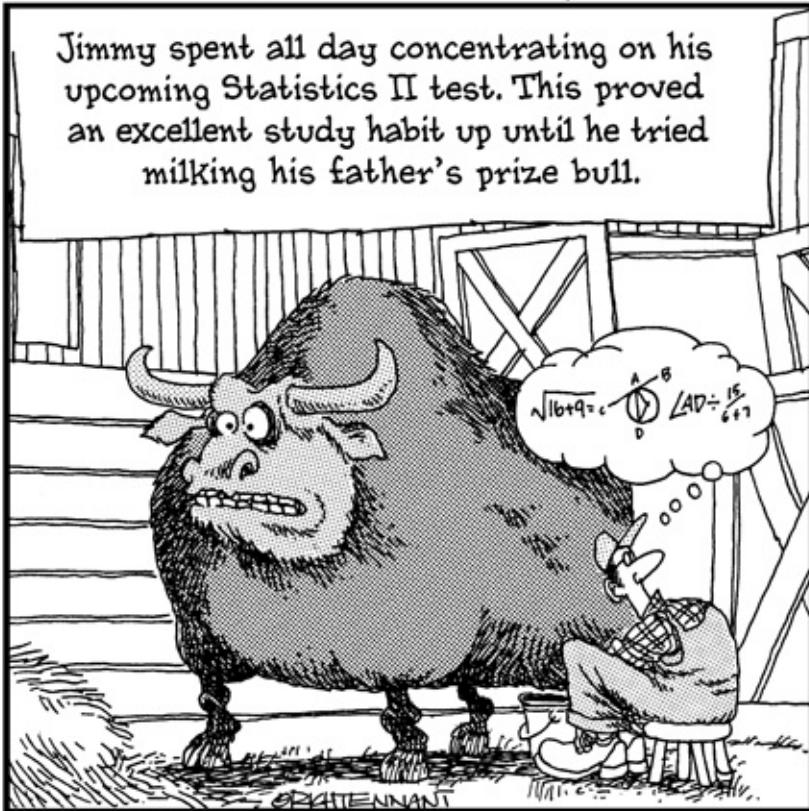
Medical scientists also think about power when they set up their studies (called *clinical trials*). Suppose they're checking to see whether an antidepressant adversely affects blood pressure (as a side effect of taking the drug). Scientists need to be able to detect small differences in blood pressure, because for some patients, any change in blood pressure is important to note and treat.

Part II

Using Different Types of Regression to Make Predictions

The 5th Wave

By Rich Tennant



In this part...

This part takes you beyond using one variable to predict another variable using a straight line (that's what simple linear regression is about). Instead, you find ways to predict one variable using many, and you also discover ways to make predictions using curves. Finally, you make predictions for probabilities, not just average values. It's a one-stop-shopping place for all things regression. Because these methods allow you to solve more complex problems, they lend themselves nicely to many real-world applications.

Chapter 4

Getting in Line with Simple Linear Regression

In This Chapter

- ▶ Using scatterplots and correlation coefficients to examine relationships
 - ▶ Building a simple linear regression model to estimate y from x
 - ▶ Testing how well the model fits
 - ▶ Interpreting the results and making good predictions
-

Looking for relationships and making predictions is one of the staples of data analysis. Everyone wants to answer questions like, “Can I predict how many units I’ll sell if I spend x amount of advertising dollars?”; or “Does drinking more diet cola really relate to more weight gain?”; or “Do children’s backpacks seem to get heavier with each year of school, or is it just me?”

Linear regression tries to find relationships between two or more variables and comes up with a model that tries to describe that relationship, much like the way the line $y = 2x + 3$ explains the relationship between x and y . But unlike in math, where functions like $y = 2x + 3$ tell the entire story about the two variables, in statistics things don’t come out that perfectly; some variability and error is involved (that’s what makes it fun!).

This chapter is partly a review of the concepts of simple linear regression presented in a typical Stats I textbook. But the fun doesn’t stop there. I expand on the ideas about regression that you picked up in your Stats I course and set you up for some of the other types of regression models you see in Chapters 5 through 8.

In this chapter, you see how to build a simple linear regression model that examines the relationship between two variables. You also see how simple linear regression works from a model-building standpoint.

Exploring Relationships with Scatterplots and Correlations

Before looking ahead to predicting a value of y from x using a line, you need to

- ✓ Establish that you have a legitimate reason to do so by using a straight line.
- ✓ Feel confident that using a line to make that prediction will actually work well.

In order to accomplish both of these important steps, you need to first plot the data in a pairwise fashion so you can visually look for a relationship; then you need to somehow

quantify that relationship in terms of how well those points follow a line. In this section, you do just that, using scatterplots and correlations.

Here's a perfect example of a situation where simple linear regression is useful: In 2004, the California State Board of Education wrote a report entitled "Textbook Weight in California: Analysis and Recommendations." This report discussed the great concern over the weight of the textbooks in students' backpacks and the problems it presents for students. The board conducted a study where it weighed a variety of textbooks from each of four core areas studied in grades 1–12 (reading, math, science, and history — where's statistics?) over a range of textbook brands and found the average total weight for all four books for each grade.

The board consulted pediatricians and chiropractors, who recommended that the weight of a student's backpack should not exceed 15 percent of his or her body weight. From there, the board hypothesized that the total weight of the textbooks in these four areas increases for each grade level and wanted to see whether it could find a relationship between the average child's weight in each grade and the average weight of his or her books. So along with the average weight of the four core-area textbooks for each grade, researchers also recorded the average weight for the students in that grade. The results are shown in Table 4-1.

Table 4-1 Average Textbook Weight and Student Weight (Grades 1–12)

Grade	Average Student Weight (In Pounds)	Average Textbook Weight (In Pounds)
1	48.50	8.00
2	54.50	9.44
3	61.25	10.08
4	69.00	11.81
5	74.50	12.28
6	85.00	13.61
7	89.00	15.13
8	99.00	15.47
9	112.00	17.36
10	123.00	18.07
11	134.00	20.79
12	142.00	16.06

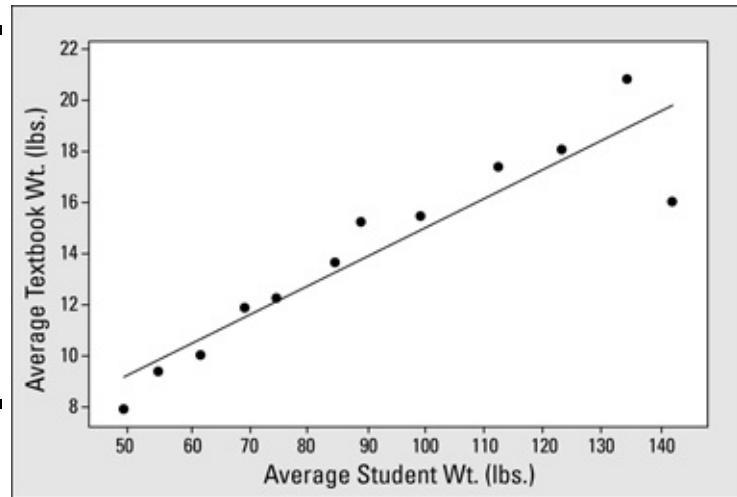
In this section, you begin exploring whether or not a relationship exists between these two quantitative variables. You start by displaying the pairs of data using a two-dimensional scatterplot to look for a possible pattern, and you quantify the strength and direction of that pattern using the correlation coefficient.

Using scatterplots to explore relationships

In order to explore a possible relationship between two variables, such as textbook weight and student weight, you first plot the data in a special graph called a *scatterplot*. A scatterplot is a two-dimensional graph that displays pairs of data, one pair per observation in the (x, y) format. Figure 4-1 shows a scatterplot of the textbook-weight data from Table 4-1.

You can see that the relationship appears to follow the straight line that's included on the graph, except possibly for the last point, where textbook weight is 16.06 pounds and student weight is 142 pounds (for grade 12). This point appears to be an *outlier* — it's the only point that doesn't fall into the pattern. So overall, an uphill, or *positive* linear relationship appears to exist between textbook weight and student weight; as student weight increases, so does textbook weight.

Figure 4-1:
Scatterplot of
average
student
weight
versus
average
textbook
weight in
grades 1–12.



To make a scatterplot in Minitab, enter the data in columns one and two of the spreadsheet. Go to Graphs>Scatterplot. Click Simple and then OK. Highlight the response variable (y) in the left-hand box, and click Select. This variable shows up as the y variable in the scatterplot. Click on the explanatory (x) variable in the left-hand box, and click Select. It shows up in the x variable box. Click OK, and you get the scatterplot.

Collating the information by using the correlation coefficient

After you've displayed the data using a scatterplot (see the preceding section), the next step is to find a statistic that quantifies the relationship somehow. The *correlation coefficient* (also known as *Pearson's correlation coefficient*, especially in statistical software packages) measures the strength and direction of the linear relationship between two quantitative variables x and y . It's a number between -1 and $+1$ that's *unit-free*, which means that if you change from pounds to ounces, the correlation coefficient doesn't change. (What a messed-up world it would be if this wasn't the case!)

If the relationship between x and y is uphill, or positive (as x increases, so does y), the correlation is a positive number. If the relationship is downhill, or negative (as x increases, y gets smaller), then the correlation is negative. The following list translates different correlation values:

- ✓ **A correlation value of zero means that you can find no linear relationship between x and y .** (It may be that a different relationship exists, such as a curve; see Chapter 7 for more on this.)
- ✓ **A correlation value of +1 or -1 indicates that the points fall in a perfect, straight line.** (Negative values indicate a downhill relationship; positive values indicate an uphill relationship.)
- ✓ **A correlation value close to +1 or -1 signifies a strong relationship.** A general rule of thumb is that correlations close to or beyond 0.7 or -0.7 are considered to be strong.
- ✓ **A correlation closer to +0.5 or -0.5 shows a moderate relationship.**

You can calculate the correlation coefficient by using a formula involving the standard deviation of x , the standard deviation of y , and the covariance of x and y , which measures how x and y move together in relation to their means. However, the formula isn't the focus here (you can find it in your Stats I textbook or in my other book, *Statistics For Dummies*, published by Wiley); it's the concept that's important. Any computer package can calculate the correlation coefficient for you with a simple click of the mouse.



To have Minitab calculate a correlation for you, go to Stat>Basic Statistics>Correlation. Highlight the variables you want correlations for, and click Select. Then click OK.

The correlation for the textbook-weight example is (can you guess before looking at it?) 0.926, which is very close to 1.0. This correlation means that a very strong linear relationship is present between average textbook weight and average student weight for grades 1–12, and that relationship is positive and linear (it follows a straight line). This correlation is confirmed by the scatterplot shown in Figure 4-1.



Data analysts should never make any conclusions about a relationship between x and y based solely on either the correlation or the scatterplot alone; the two elements need to be examined together. It's possible (but of course not a good idea) to manipulate graphs to look better or worse than they really are just by changing the scales on the axes. Because of this, statisticians never go with the scatterplot alone to determine whether or not a linear relationship exists between x and y . A correlation without a scatterplot is dangerous, too, because the relationship

between x and y may be very strong but just not linear.

Building a Simple Linear Regression Model

After you have a handle on which x variables may be related to y in a linear way, you go about the business of finding that straight line that best fits the data. You find the slope and y -intercept, put them together to make a line, and you use the equation of that line to make predictions for y . All this is part of building a simple linear regression model.

In this section, you set the foundation for regression models in general (including those you can find in Chapters 5 through 8). You plot the data, come up with a model that you think makes sense, assess how well it fits, and use it to guesstimate the value of y given another variable, x .

Finding the best-fitting line to model your data

After you've established that x and y have a strong linear relationship, as evidenced by both the scatterplot and the correlation coefficient (close to or beyond 0.7 and -0.7 ; see the previous sections), you're ready to build a model that estimates y using x . In the textbook-weight case, you want to estimate average textbook weight using average student weight.

The most basic of all the regression models is the simple linear regression model that comes in the general form of $y = \alpha + \beta x + \varepsilon$. Here, α represents the y -intercept of the line, β represents the slope, and ε represents the error in the model due to chance.



A straight line that's used in simple linear regression is just one of an entire family of models (or functions) that statisticians use to express relationships between variables. A *model* is just a general name for a function that you can use to describe what outcome will occur based on some given information about one or more related variables.

Note that you will never know the true model that describes the relationship perfectly. The best you can do is estimate it based on data.

To find the right model for your data, the idea is to scour all possible lines and choose the one that fits the data best. Thankfully, you have an algorithm that does this for you (computers use it in their calculations). Formulas also exist for finding the slope and y -intercept of the best-fitting line by hand. The best-fitting line based on your data is $y = a + bx$, where a estimates α and b estimates β from the true model. (You can find those formulas in your Stats I text or in *Statistics For Dummies*.)



To run a linear regression analysis in Minitab, go to Stat>Regression>Regression. Highlight the response (y) variable in the left-hand box, and click Select. The variable shows up in the Response Variable box. Then highlight your explanatory (x) variable, and click Select. This variable shows up in the Predictor Variable box. Click OK.

The equation of the line that best describes the relationship between average textbook weight and average student weight is $y = 3.69 + 0.113x$, where x is the average student weight for that grade, and y is the average textbook weight. Figure 4-2 shows the Minitab output of this analysis.

Figure 4-2:
Simple linear regression analysis for the textbook-weight example.

The regression equation is textbook wt = 3.69 + 0.113 student wt					
 Predictor Coef SE Coef T P					
Constant	3.694	1.395	2.65	0.024	
student wt	0.11337	0.01456	7.78	0.000	
 $S = 1.51341$ R-Sq = 85.8% R-Sq(adj) = 84.4%					



By writing $y = 3.69 + 0.113x$, you mean that this equation represents your estimated value of y , given the value of x that you observe with your data.

Statisticians technically write this equation by using a caret (or *hat* as statisticians call it), like \hat{y} , so everyone can know it's an estimate, not the actual value of y . This y -hat is your estimate of the average value of y over the long term, based on the observed values of x . However, in many Stats I texts, the hat is left off because statisticians have an unwritten understanding as to what y represents. This issue comes up again in Chapters 5 through 8. (By the way, if you think y -hat is a funny term here, it's even funnier in Mexico, where statisticians call it *y-sombrero* — no kidding!)

The y -intercept of the regression line

Selected parts of that Minitab output shown in Figure 4-2 are of importance to you at this point. First, you can see that under the Coef column you have the numerical values on the right side of the equation of the line — in other words, the slope and y -intercept. The number 3.69 represents the coefficient of “Constant,” which is a fancy way of saying that’s the y -intercept (because the y -intercept is just a constant — it never changes). The y -intercept is the point where the line crosses the y -axis; in other words, it’s the value of y when x equals zero.



The y -intercept of a regression line may or may not have a practical meaning depending on the situation. To determine whether the y -intercept of a regression line has practical meaning, look at the following:

- ✓ Does the y -intercept fall within the actual values in the data set? If yes, it has practical meaning.
- ✓ Does the y -intercept fall into negative territory where negative y -values aren't possible? For example, if the y -values are weights, they can't be negative. Then the y -intercept has no practical meaning. The y -intercept is still needed in the equation though, because it just happens to be the place where the line, if extended to the y -axis, crosses the y -axis.
- ✓ Does the value $x = 0$ have practical meaning? For example, if x is temperature at a football game in Green Bay, then $x = 0$ is a value that's relevant to examine. If $x = 0$ has practical meaning, then the y -intercept does too, because it represents the value of y when $x = 0$. If the value of $x = 0$ doesn't have practical meaning in its own right (such as when x represents height of a toddler), then the y -intercept doesn't either.

In the textbook example, the y -intercept doesn't really have a practical meaning because students don't weigh zero pounds, so you don't really care what the estimated textbook weight is for that situation. But you do need to find a line that fits the data you do have (where average student weights go from 48.5 to 142 pounds). That best-fitting line must include a y -intercept, and for this problem, that y -intercept happens to be 3.69 pounds.

The slope of the regression line

The value 0.113 from Figure 4-2 indicates the coefficient (or number in front) of the student-weight variable. This number is also known as the *slope*. It represents the change in y (textbook weight) is associated with a one-unit increase in x (student weight). As student weight increases by 1 pound, textbook weight increases by about 0.113 pounds, on average. To make this relationship more meaningful, you can multiply both quantities by 10 to say that as student weight increases by 10 pounds, the textbook weight goes up by about 1.13 pounds on average.



Whenever you get a number for the slope, take that number and put it over 1 to help you get started on a proper interpretation of slope. For example, a slope of 0.113 is rewritten as $\frac{0.113}{1}$. Using the idea that slope equals rise over run, or change in y over change in x , you can interpret the value of 0.113 in the following way: As x increases on average by 1 pound, y increases by 0.113 pounds.

Making point estimates by using the regression line

When you have a line that estimates y given x , you can use it to give a one-number estimate for the (average) value of y for a given value of x . This is called making a *point estimate*. The basic idea is to take a reasonable value of x , plug it into the equation of the regression line, and see what you get for the value of y .

In the textbook-weight example, the best-fitting line (or model) is the line $y = 3.69 + 0.113x$. For an average student who weighs 60 pounds, for example, a one-number point estimate of the average textbook weight is $3.69 + (0.113 * 60) = 10.47$ pounds (those poor little kids!). If the average student weighs 100 pounds, the estimated average textbook weight is $3.69 + (0.113 * 100) = 14.99$, or nearly 15 pounds, plus or minus something. (You find out what that something is in the following section.)

No Conclusion Left Behind: Tests and Confidence Intervals for Regression

After you have the slope of the best-fitting regression line for your data (see the previous sections), you need to step back and take into account the fact that sample results will vary. You shouldn't just say, "Okay, the slope of this line is 2. I'm done!" It won't be exactly 2 the next time. This variability is why statistics professors harp on adding a margin of error to your sample results; you want to be sure to cover yourself by adding that plus or minus.

In hypothesis testing, you don't just compare your sample mean to the population mean and say, "Yep, they're different alright!" You have to standardize your sample result using the standard error so that you can put your results in the proper perspective (see Chapter 3 for a review of confidence intervals and hypothesis tests).

The same idea applies here with regression. The data were used to figure out the best-fitting line, and you know it fits well for that data. That's not to say that the best-fitting line will work perfectly well for a new data set taken from the same population. So, in regression, all your results should involve the standard error with them in order to allow for the fact that sample results vary. That goes for estimating and testing for the slope and y -intercept and for any predictions that you make.

Many times in Stats I courses the concept of margin of error is skipped over after the best-fitting regression line is found. But these are very important ideas and should always be included. (Okay, enough of the soap box for now. Let's get out there and do it!)

Scrutinizing the slope

Recall the *slope* of the regression line is the amount by which you expect the y variable to change on average as the x variable increases by 1 unit — the old rise-over-run idea (see the section "The slope of the regression line" earlier in this chapter). Now, how do you deal with knowing the best-fitting line will change with a new data set? You just apply the

basic ideas of confidence intervals and hypothesis tests (see Chapter 3).

A confidence interval for slope

A *confidence interval* in general has this form: your statistic plus or minus a margin of error. The margin of error includes a certain number of standard deviations (or standard errors) from your statistic. How many standard errors you add and subtract depends on what confidence level, $1 - \alpha$, you want. The size of the standard error depends on the sample size and other factors.

The equation of the best-fitting simple linear regression line, $y = a + bx$, includes a slope (b) and a y -intercept (a). Because these were found using the data, they're only estimates of what's going on in the population, and therefore they need to be accompanied by a margin of error.

The formula for a $1 - \alpha$ level confidence interval for the slope of a regression line is

$$b \pm t_{n-2}^* * SE_b, \text{ where the standard error is denoted } SE_b = \frac{s}{\sqrt{\sum_i (x_i - \bar{x})^2}}, \text{ where } s = \sqrt{\frac{1}{n-2} \sum_i (y_i - \hat{y}_i)^2}$$

. The value of t^* comes from the t -distribution with $n - 2$ degrees of freedom and area to its right equal to $\alpha \div 2$. (See Chapter 3 regarding the concept of α .)



In case you wonder why you see $n - 2$ degrees of freedom here, as opposed to $n - 1$ degrees of freedom used in t -tests for the population mean in Stats I, here's the scoop. From Stats I you know that a *parameter* is a number that describes the population; it's usually known, and it can change from scenario to scenario. For each parameter in a model, you lose 1 degree of freedom. The regression line contains two parameters — the slope and the y -intercept — and you lose 1 degree of freedom for each one. With the t -test from Stats I you only have one parameter, the population mean, to worry about, hence you use $n - 1$ degrees of freedom.



You can find the value of t^* in any t -distribution table (check your textbook for one). For example, suppose you want to find a 95 percent confidence interval based on sample size $n = 10$. The value of t^* is found in Table A-1 in the appendix in the row marked $10 - 2 = 8$ degrees of freedom, and the column marked 0.025 (because $\alpha \div 2 = 0.05 \div 2 = 0.025$). This value of t^* is 2.306. (*Statistics For Dummies* can tell you a lot more about the t -distribution and the t -table.)

To put together a 95 percent confidence interval for the slope using computer output, you pull off the pieces that you need. For the textbook-weight example, in Figure 4-2 you see that the slope is equal to 0.11337. (Recall that slope is the coefficient of the x variable in the equation, which is why you see the abbreviation Coef in the output.)

Because the slope changes from sample to sample, it's a random variable with its own distribution, its own mean, and its own standard error. (Recall from Stats I the standard error of a statistic is likened to the standard deviation of a random variable.) If you look just to the right of the slope in Figure 4-2, you see SE Coef; this stands for the standard error of the slope (which is 0.01456 in this case).

Now all you need is the value of t^* from the t -table (Table A-1 in the appendix). Because $n = 12$, you look in the row where degrees of freedom is $12 - 2 = 10$. You want a 95 percent confidence interval, so you look in the column for $(1 - 0.95) \div 2 = 0.25$. The t^* value you get is 2.228.

Putting these pieces together, a 95 percent confidence interval for the slope of the best-fitting regression line for the textbook-weight example is $0.11337 \pm 2.228 * 0.01456$ which goes from 0.0809 to 0.1458. The units are in pounds (textbook) per pounds (child weight). Note this interval is large due to the small sample size, which increases the standard error.

A hypothesis test for slope

You may be interested in conducting a hypothesis test for the slope of a regression line as another way to assess how well the line fits. If the slope is zero or close to it, the regression line is basically flat, signifying that no matter the value of x , you'll always estimate y by using its mean. This means that x and y aren't related at all, so a specific value of x doesn't help you predict a specific value for y . You can also test to see if the slope is some value other than zero, but that's atypical. So for all intents and purposes, I use the hypotheses $H_0: \beta = 0$ versus $H_a: \beta \neq 0$, where β is the slope of the true model.

To conduct a hypothesis test for the slope of a simple linear regression line, you follow the basic steps of any hypothesis test. You take the statistic (b) from your data, subtract the value in H_0 (in this case it's zero), and standardize it by dividing by the standard error (see Chapter 3 for more on this process).

Using the formula for standard error for b , the test statistic for the hypothesis test of

whether or not the slope equals zero is $\frac{b-0}{SE_b}$, where $SE_b = \frac{s}{\sqrt{\sum_i (x_i - \bar{x})^2}}$, and $s = \sqrt{\frac{1}{n-2} \sum_i (y_i - \hat{y}_i)^2}$. On the Minitab output from Figure 4-2, the test statistic is located right next to the SE Coef column; it's cleverly marked T. In this case $T = 7.78$. Compare this value to $t^* = 2.228$ from the t-table. Because $T > t^*$, you have strong evidence to reject H_0 and conclude that the slope of the regression line for the textbook-weight data is not zero. (In fact, it has to be greater than that, according to your data.)

You can also just find the exact p -value on the output, right next to the T column, in the column is marked P. The p -value for the test for slope in this case is 0.000, which means it's less than 0.001. You conclude that the slope of this line is not zero, so textbook weight is significantly related to student weight. (See Chapter 3 to brush up on p -values.)



To test to see whether the slope is some value other than zero, just plug that value in for b_0 in the formula for the test statistic. Also, you may conduct one-sided hypothesis tests to see whether the slope is strictly greater than zero or strictly less than zero. In those cases, you find the same test statistic but compare it to the value t^* where the area to the right (or left, respectively) is α .

Inspecting the y-intercept

The y-intercept is the place where the regression line $y = a + bx$ crosses the y -axis and is denoted by a (see the earlier section “The y-intercept of the regression line”). Sometimes the y-intercept can be interpreted in a meaningful way, and sometimes not. This differs from slope, which is always interpretable. In fact, between the two elements of slope and intercept, the slope is the star of the show, with the y-intercept serving as the somewhat less famous but still noticeable sidekick.

There are times when the y-intercept makes no sense. For example, suppose you use rain to predict bushels per acre of corn; if you have zero rain, you have zero corn, but if the regression line crosses the y -axis somewhere else besides zero (and it most likely will), the y-intercept will make no sense. Another situation is where no data were collected near the value of $x = 0$; interpreting the y-intercept at that point is not appropriate. For example, using a student’s score on midterm 1 to predict her score on midterm 2, unless the student didn’t take the exam at all (in which case it doesn’t count), she’ll get at least some points.

Many times, however, the y-intercept is of interest to you and has a value that you can interpret, such as when you’re talking about predicting coffee sales using temperature for football games. Some games get cold enough to have zero and subzero temperatures (like Packers games for example — Go Pack Go!).

Suppose I collect data on ten of my students who recorded their study time (in minutes) for a 10-point quiz, along with their quiz scores. The data have a strong linear relationship by all the methods used in this chapter (for example, refer to the earlier section “Exploring Relationships with Scatterplots and Correlations”). I went ahead and conducted a regression analysis, and the results are shown in Figure 4-3.

Because there are students who (heaven forbid!) didn’t study at all for the quiz, the y-intercept of 3.29 points (where study time $x = 0$) can be interpreted safely. Its value is shown in the Coef column in the row marked Constant (see the section “The y-intercept of the regression line” for more information). The next step is to give a confidence interval for the y-intercept of the regression line, where you can take conclusions beyond just this sample of ten students.

The formula for a $1 - \alpha$ level confidence interval for the y-intercept (a) of a simple linear

regression line is $a \pm t^*_{n-2} SE_a$. The standard error, SE_a , is equal to

$$SE_a = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_i (x_i - \bar{x})^2}}, \text{ where } s = \sqrt{\frac{1}{n-2} \sum_i (y_i - \hat{y}_i)^2}$$

, where again the value of t^* comes from the t -

distribution with $n - 2$ degrees of freedom whose area to the right is equal to $\alpha / 2$. Using the output from Figure 4-3 and the t -table, I'm 95 percent confident that the quiz score (y) for someone with a study time of $x = 0$ minutes is $3.29 \pm 2.306 * 0.4864$, which is anywhere from 2.17 to 4.41, on average. Note that 2.306 comes from the t -table with $10 - 2 = 8$ degrees of freedom and 0.4864 is the SE for the y -intercept from Figure 4-3. (So studying for zero minutes for my quiz is not something to aspire to.)

By the way, to find out how much time studying affected the quiz score for these students, you can get an estimate of the slope on the output from Figure 4-3 that the coefficient for slope is 0.1793, which says each minute of studying is related to an increase in score of 0.1793 of a point, plus or minus the margin of error, of course. Or, 10 more minutes relates to 1.793 more points. On a 10-point quiz, it all adds up!

Figure 4-3:
Regression analysis for study time and quiz score data.

The regression equation is quiz score = 3.29 + 0.179 minutes studying					
Predictor	Coef	SE Coef	T	P	
Constant	3.2931	0.4864	6.77	0.000	
Minutes studying	0.17931	0.02103	8.53	0.000	
$S = 0.877153$		$R-Sq = 90.1\%$		$R-Sq (adj) = 88.8\%$	



Testing a hypothesis about the y -intercept isn't really something you'll find yourself doing much because most of the time you don't have a preconceived notion about what the y -intercept would be (nor do you really care ahead of time). The confidence interval is much more useful. However, if you do need to conduct a hypothesis test for the y -intercept, you take your y -intercept, subtract the value in H_0 , and divide by the standard error, found on the computer output in the row for Constant and the column for SE Coef. (The default value is to test to see whether the y -intercept is zero.) The test is in the T column of the output, and its p -value is shown in the P column. In the study time and quiz score example, the p -value is 0.000, so the y -intercept is significantly different from zero. All this means is that the line crosses the y -axis somewhere else.

Building confidence intervals for the average response

When you have the slope and y -intercept for the best-fitting regression line, you put them together to get the line $y = a + bx$. The value of y here really represents the average value of y for a particular value of x . For example, in the textbook-weight data, Figure 4-2 shows the regression line $y = 3.69 + 0.11337x$ where $x = \text{average student weight}$ and $y = \text{average}$

textbook weight. If you put in 100 pounds for x , you get $y = 3.69 + 0.1137 * 100 = 15.02$ pounds of textbook weight for the group averaging 100 pounds. This number, 15.02, is an estimate of the average weight of textbooks for children of this weight.

But you can't stop there. Because you're getting an estimate of the average textbook weight using y , you also need a margin of error for y to go with it, to create a confidence interval for the average y at a given x that generalizes to the population.

Take your estimate, y , which you get by plugging your given x value into the regression line, and then add and subtract the margin of error for y . The formula for a $1 - \alpha$ confidence interval for the mean of y for a given value of x (call it x^*) is equal to $y \pm t_{n-2}^* SE_{\hat{\mu}}$, where y is the value of the equation of the line when you plug in x^* for x . The standard

$$SE_y = s \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_i (x_i - \bar{x})^2}}, \text{ where } s = \sqrt{\frac{1}{n-2} \sum_i (y_i - \hat{y}_i)^2}$$

error for y is equal to . Luckily Minitab does these calculations for you and reports a confidence interval for the mean of y for a given x^* .



To find a confidence interval for the mean value of y using Minitab, you ask for a regression analysis (see instructions in the earlier section "Finding the best-fitting line to model your data") and click on Options. You see a box called Prediction Intervals for New Observations; enter the value of x^* that you want, and just below that, put in your confidence level (the default is 95 percent). Even though this box is labeled Prediction Intervals, it finds both a confidence interval and a prediction interval for you as well. (Prediction intervals are different from confidence intervals, and I discuss them in the next section.) On the computer output, the confidence interval is labeled 95% CI.

Returning to the textbook-weight example, the computer output for finding a 95 percent confidence interval for the average textbook weight for 100-pound children is shown in Figure 4-4. The result is (14.015, 16.048) pounds. I'm 95 percent confident that the average textbook weight for the group of children averaging 100 pounds is between 14.015 and 16.048 pounds. (Get out those rolling backpacks, kids!)



You should only make predictions for the average value of y for x values that are within the range of where the data was collected. Failure to do so will result in the statistical no-no called extrapolation (see the later section "Knowing the Limitations of Your Regression Analysis").

Making the band with prediction intervals

Suppose instead of the mean value of y , you want to take a guess at what y would be for some future value of x . Because you're looking into the future, you have to make a

prediction, and to do that, you need a range of likely values of y for a given x^* . This is what statisticians call a *prediction interval*.

The formula for a $1 - \alpha$ level prediction interval for y at a given value x^* is

$$SE_y = s \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_i (x_i - \bar{x})^2}}$$

, where $s = \sqrt{\frac{1}{n-2} \sum_i (y_i - \hat{y}_i)^2}$

for you.



To find a $1 - \alpha$ level prediction interval for the value of y for a given x^* using Minitab, you ask for a regression analysis (see instructions in the earlier section “Finding the best-fitting line to model your data”) and click Options. In the box Prediction Intervals for New Observations, enter the value of x that you want, and just below that, put in your confidence level (the default is 95 percent). On the computer output, the prediction interval is labeled 95% PI, and it appears right next to the confidence interval for the mean of y for that same x^* .

Predicting textbook weight using student weight

For the textbook-weight data, suppose you’ve already made your regression line and now a new student comes on the scene. You want to predict this student’s textbook weight. This means you want a prediction interval rather than a confidence interval, because you want to predict the textbook weight for one person, not the average weight for a group.

Suppose this new student weighs 100 pounds. To find the prediction interval for the textbook weight for this student, you use $x^* = 100$ pounds and let Minitab do its thing.

The computer output in Figure 4-4 shows the 95 percent prediction interval for textbook weight for a single 100-pound child is (11.509, 18.533) pounds. Note this is wider than the confidence interval of (14.015, 16.048) for the mean textbook weight for 100-pound children found in the earlier section “Building confidence intervals for the average response.” This difference is due to the increased variability in looking at one child and predicting one textbook weight.

Figure 4-4:
Prediction
interval of
textbook
weight for a
100-pound
child.

Predicted Values for New Observations				
New				
Obs	Fit	SE Fit	95% CI	95% PI
1	15.031	0.456	(14.015, 16.048)	(11.509, 18.533)

Comparing prediction and confidence intervals

Note that the formulas for prediction intervals and confidence intervals are very similar.

In fact, the prediction interval formula is exactly the same as the confidence interval formula except it adds a 1 under the square root. Because of this difference in the formulas, the margin of error for a prediction interval is larger than for a confidence interval.

This difference also makes sense from a statistical point. A prediction interval has more variability than a confidence interval because it's harder to make a prediction about y for a single value of x^* than it is to estimate the average value of y for a given x^* . (For example, individual test scores vary more than average test scores do.) A prediction interval will be wider than a confidence interval; it will have a larger margin of error.

A similarity between prediction intervals and confidence intervals is that their margin of error formulas both contain x^* , which means the margin of error in either case depends on which value of x^* you use. It turns out in both cases that if you use the mean value of x as your x^* , the margin of error for each interval is at its smallest because there's more data around the mean of x than at any other value. As you move away from the mean of x , the margin of error increases for each interval.

Checking the Model's Fit (The Data, Not the Clothes!)

After you've established a relationship between x and y and have come up with an equation of a line that represents that relationship, you may think your job is done. (Many researchers erringly stop here, so I'm depending on you to break the cycle!) The most-important job remains to be completed: checking to be sure that the conditions of the model are truly met and that the model fits well in more specific ways than the scatterplot and correlation measure (which I cover in the earlier section "Exploring Relationships with Scatterplots and Correlations").

This section presents methods for defining and assessing the fit of a simple linear regression model.

Defining the conditions

Two major conditions must be met before you apply a simple linear regression model to a data set:

- ✓ The y 's must have an approximately normal distribution for each value of x .
- ✓ The y 's must have a constant amount of spread (standard deviation) for each value of x .

Normal y's for every x

For any value of x , the population of possible y -values must have a normal distribution.

The mean of this distribution is the value for y that's on the best-fitting line for that x -value. That is, some of your data fall above the best-fitting line, some data fall below the best fitting line, and a few may actually land right on the line.

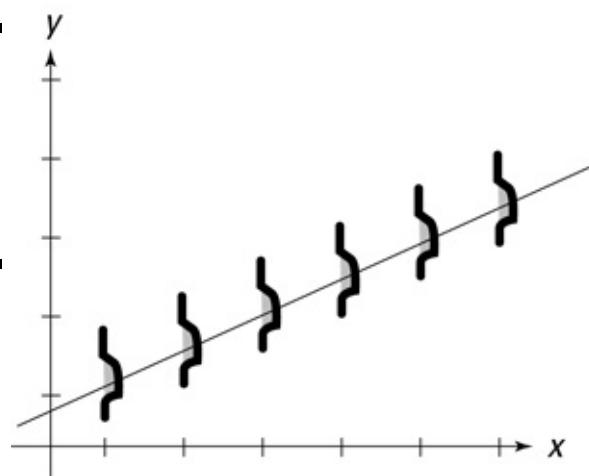


If the regression model is fitting well, the data values should be scattered around the best-fitting line in such a way that about 68 percent of the values lie within one standard deviation of the line, about 95 percent of the values lie within two standard deviations of the line, and about 99.7 percent of the values lie within three standard deviations of the line. This specification, as you may recall from your Stats I course, is called the *68-95-99.7 rule*, and it applies to all bell-shaped data (for which the normal distribution applies).

You can see in Figure 4-5 how for each x -value, the y -values you may observe tend to be located near the best-fitting line in greater numbers, and as you move away from the line, you see fewer and fewer y -values, both above and below the line. More than that, they're scattered around the line in a way that reflects a bell-shaped curve, the normal distribution. This indicates a good fit.

Why does this condition makes sense? The data you collect on y for any particular x -value vary from individual to individual; for example, not all students' textbooks weigh the same, even for students who weigh the exact same amount. But those values aren't allowed to vary any way they want to. To fit the conditions of a linear regression model, for each given value of x , the data should be scattered around the line according to a normal distribution. Most of the points should be close to the line, and as you get farther from the line, you can expect fewer data points to occur. So condition number one is that the data have a normal distribution for each value of x .

Figure 4-5:
Conditions of
a simple
linear
regression
model.



Same spread for every x

In order to use the simple linear regression model, as you move from left to right on the x -axis, the spread in the y -values around the line should be the same, no matter which value of x you're looking at. This requirement is called the *homoscedasticity condition*.

(How they came up with that mouthful of a word just for describing the fact that the standard deviations stay the same across the x -values, I'll never know.) This condition ensures that the best-fitting line works well for all relevant values of x , not just in certain areas.

You can see in Figure 4-5 that no matter what the value of x is, the spread in the y -values stays the same throughout. If the spread got bigger and bigger as x got larger and larger, for example, the line would lose its ability to fit well for those large values of x .

Finding and exploring the residuals

To check to see whether the y -values come from a normal distribution, you need to measure how far off your predictions were from the actual data that came in. These differences are called *errors*, or *residuals*. To evaluate whether a model fits well, you need to check those errors and see how they stack up.



In a model-fitting context, the word error doesn't mean "mistake." It just means a difference between the data and the prediction based on the model. The word I like best to describe this difference is residual, however. It sounds more upbeat.

The following sections focus on finding a way to measure these residuals that the model makes. You also explore the residuals to identify particular problems that occurred in the process of trying to fit a straight line to the data. In other words, you can discover that looking at residuals helps you assess the fit of the model and diagnose problems that caused a bad fit, if that was the case.

Finding the residuals

A *residual* is the difference between the observed value \hat{y} of y (from the best-fitting line) and the predicted value of y , *also known as* y (from the data set). Its notation is $(y - \hat{y})$. Specifically, for any data point, you take its observed y -value (from the data) and subtract its expected y -value (from the line). If the residual is large, the line doesn't fit well in that spot. If the residual is small, the line fits well in that spot.

For example, suppose you have a point in your data set $(2, 4)$ and the equation of the best-fitting line is $y = 2x + 1$. The expected value of y in this case is $(2 * 2) + 1 = 5$. The observed value of y from the data set is 4. Taking the observed value minus the estimated value, you get $4 - 5 = -1$. The residual for that particular data point $(2, 4)$ is -1 . If you observe a y -value of 6 and use the same straight line to estimate y , then the residual is $6 - 5 = +1$.



In general, a positive residual means you underestimated y at that point; the line

is below the data. A negative residual means you overestimated y at that point; the line is above the data.

Standardizing the residuals

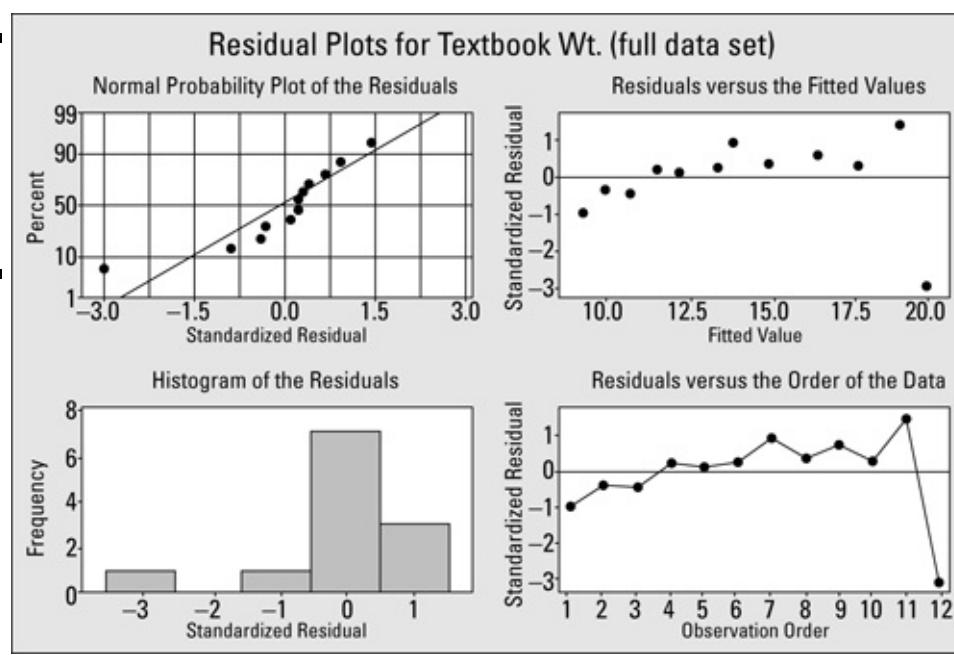
Residuals in their raw form are in the same units as the original data, making them hard to judge out of context. To make interpreting the residuals easier, statisticians typically *standardize* them — that is, subtract the mean of the residuals (zero) and divide by the standard deviation of all the residuals. The residuals are a data set just like any other data set, so you can find their mean and standard deviation like you always do.

Standardizing just means converting to a Z-score so that you see where it falls on the standard normal distribution. (See your Stats I text or *Statistics For Dummies* for info on Z-scores.)

Making residual plots

You can plot the residuals on a graph called a *residual plot*. (If you've standardized the residuals, you call it a *standardized residual plot*.) Figure 4-6 shows the Minitab output for a variety of standardized residual plots, all getting at the same idea: checking to be sure the conditions of the simple linear regression model are met.

Figure 4-6:
Standardized residual plots for textbook-weight data.



Checking normality

If the condition of normality is met, you can see on the residual plot lots of (standardized) residuals close to zero; as you move farther away from zero, you can see fewer residuals. **Note:** You shouldn't expect to see a standardized residual at or beyond $+3$ or -3 . If this occurs, you can consider that point an outlier, which warrants further investigation. (For more on outliers, see the section “Scoping for outliers” later in this chapter.)



The residuals should also occur at random — some above the line, and some below the line. If a pattern occurs in the residuals, the line may not be fitting right.

The plots in Figure 4-6 seem to have an issue with the very last observation, the one for 12th graders. In this observation, the average student weight (142) seemed to follow the pattern of increasing with each grade level, but the textbook weight (16.06) was less than for 11th graders (20.79) and is the first point to break the pattern.

You can also see in the plot in the upper-right corner of Figure 4-6 that the very last data value has a standardized residual that sticks out from the others and has a value of -3 (something that should be a very rare occurrence). So the value you expected for y based on your line was off by a factor of 3 standard deviations. And because this residual is negative, what you observed for y was much lower than you may have expected it to be using the regression line.

The other residuals seem to fall in line with a normal distribution, as you can see in the upper-right plot of Figure 4-6. The residuals concentrate around zero, with fewer appearing as you move farther away from zero. You can also see this pattern in the upper-left plot of Figure 4-6, which shows how close to normal the residuals are. The line in this graph represents the equal-to-normal line. If the residuals follow close to the line, then normality is okay. If not, you have problems (in a statistical sense, of course). You can see the residual with the highest magnitude is -3, and that number falls outside the line quite a bit.

The lower-left plot in Figure 4-6 makes a histogram of the standardized residuals, and you can see it doesn't look much like a bell-shaped distribution. It doesn't even look *symmetric* (the same on each side when you cut it down the middle). The problem again seems to be the residual of -3, which skews the histogram to the left.

The lower-right plot of Figure 4-6 plots the residuals in the order presented in the data set in Table 4-1. Because the data was ordered already, the lower-right residual plot looks like the upper-right residual plot in Figure 4-6, except the dots are connected. This lower-right residual plot makes the residual of -3 stand out even more.

Checking the spread of the y 's for each x

The graph in the upper-right corner of Figure 4-6 also addresses the homoscedasticity condition. If the condition is met, then the residuals for every x -value have about the same spread. If you cut a vertical line down through each x -value, the residuals have about the same spread (standard deviation) each time, except for the last x -value, which again represents grade 12. That means the condition of equal spread in the y -values is met for the textbook-weight example.



If you look at only one residual plot, choose the one in the upper-right corner of Figure 4-6, the plot of the fitted values (the values of y on the line) versus the standardized residuals. Most problems with model fit will show up on that plot because a residual is defined as the difference between the observed value of y and the fitted value of y . In a perfect world, all the fitted values have no residual at all; a large residual (such as the one where the estimated textbook weight is 20 pounds for students averaging 142 pounds; see Figure 4-1) is indicated by a point far off from zero. This graph also shows you deviations from the overall pattern of the line; for example, if large residuals are on the extremes of this graph (very low or very high fitted values), the line isn't fitting in those areas. On balance, you can say this line fits well at least for grades 1 through 11.

Using r^2 to measure model fit

One important way to assess how well the model fits is to use a statistic called the *coefficient of determination*, or r^2 . This statistic takes the value of the correlation, r , and squares it to give you a percentage. You interpret r^2 as the percentage of variability in the y variable that's explained by, or due to, its relationship with the x variable.

The y -values of the data you collect have a great deal of variability in and of themselves. You look for another variable (x) that helps you explain that variability in the y -values. After you put that x variable into the model and find that it's highly correlated with y , you want to find out how well this model did at explaining why the values of y are different.



Note that you have to interpret r^2 using different standards than those for interpreting r . Because squaring a number between -1 and $+1$ results in a smaller number (except for $+1$, -1 , and 0 , which stay the same or switch signs), an r^2 of 0.49 isn't too bad, because it's the square of $r = 0.7$, which is a fairly strong correlation.



The following are some general guidelines for interpreting the value of r^2 :

- ✓ If the model containing x explains a lot of the variability in the y -values, then r^2 is high (in the 80 to 90 percent range is considered to be extremely high). Values like 0.70 are still considered fairly high. A high percentage of variability means that the line fits well because there's not much left to explain about the value of y other than using x and its relationship to y . So a larger value of r^2 is a good thing.
- ✓ If the model containing x doesn't help much in explaining the difference in the y -values, then the value of r^2 is small (closer to zero; between, say, 0.00 and 0.30).

roughly). The model, in this case, wouldn't fit well. You need another variable to explain y other than the one you already tried.

- Values of r^2 that fall in the middle (between, say, 0.30 and 0.70) mean that x does help somewhat in explaining y , but it doesn't do the job well enough on its own. In this case, statisticians would try to add one or more variables to the model to help explain y more fully as a group (read more about this in Chapter 5).

For the textbook-weight example, the value of r (the correlation coefficient) is 0.93. Squaring this result, you get $r^2 = 0.8649$. That number means approximately 86 percent of the variability you find in average textbook weights for all students (y -values) is explained by the average student weight (x -values). This percentage tells you that the model of using year in school to estimate backpack weight is a good bet.

In the case of simple linear regression, you have only one x variable, but in Chapter 5, you can see models that contain more than one x variable. In that situation, you use r^2 to help sort out the contribution that those x variables as a group bring to the model.

Scoping for outliers

Sometimes life isn't perfect (oh really?), and you may find a residual in your otherwise tidy data set that totally sticks out. It's called an *outlier*, and it has a standardized value at or beyond +3 or -3. It threatens to blow the conditions of your regression model and send you crying to your professor.

Before you panic, the best thing to do is to examine that outlier more closely. First, can you find an error in that data value? Did someone report her age as 642, for instance? (After all, mistakes do happen.) If you do find a certifiable error in your data set, you remove that data point (or fix it if possible) and analyze the data without it. However, if you can't explain away the problem by finding a mistake, you must think of another approach.

If you can't find a mistake that caused the outlier, you don't necessarily have to trash your model; after all, it's only one data point. Analyze the data with that data point, and analyze the data again without it. Then report and compare both analyses. This comparison gives you a sense of how influential that one data point is, and it may lead other researchers to conduct more research to zoom in on the issue you brought to the surface.

In Figure 4-1, you can see the scatterplot of the full data set for the textbook-weight example. Figure 4-7 shows the scatterplot for the data set minus the outlier. The scatterplot fits the data better without the outlier. The correlation increases to 0.993, and the value of r^2 increases to 0.986. The equation for the regression line for this data set is $y = 1.78 + 0.139x$.

The slope of the regression line doesn't change much by removing the outlier (compare

it to Figure 4-2, where the slope is 0.113). However, the y -intercept changes: It's now 1.78 without the outlier compared to 3.69 with the outlier. The slopes of the lines are about the same, but the lines cross the y -axis in different places. It appears that the outlier (the last point in the data set) has quite an effect on the best-fitting line.

Figure 4-7:
Scatterplot of
textbook-
weight data
minus the
outlier.

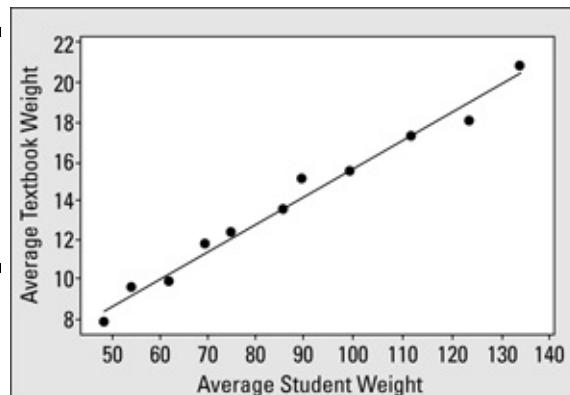
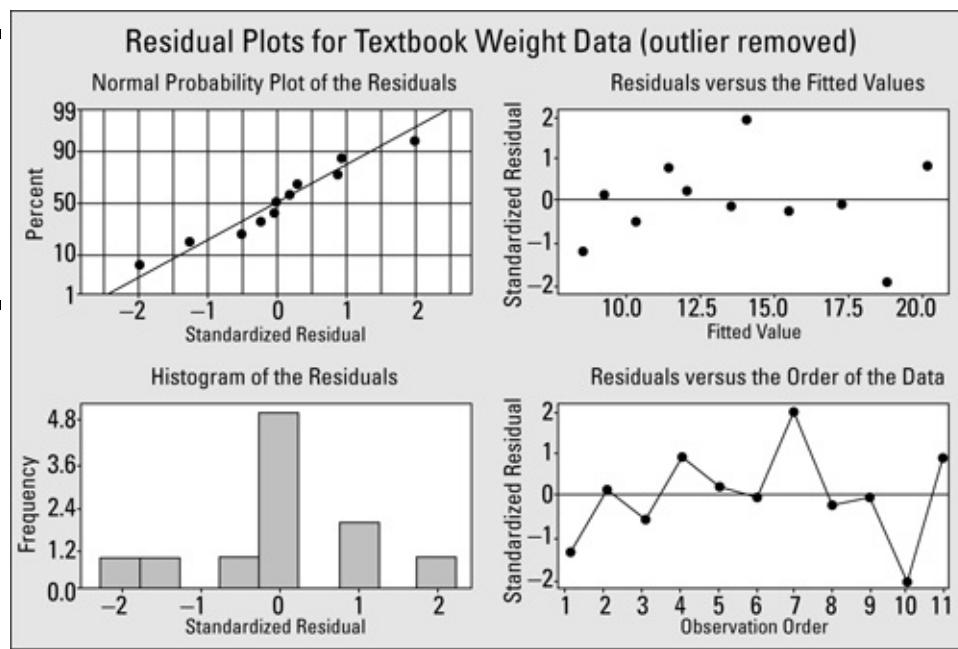


Figure 4-8 shows the residual plots for the regression line for the data set without the outlier. Each of these plots shows a much better fit of the data to the model compared to Figure 4-6. This result tells you that the data for grade 12 is influential in this data set and that the outlier needs to be noted and perhaps explored further. Do students peak when they're juniors in high school? Or do they just decide when they're seniors that it isn't cool to carry books around? (A statistician's job isn't to wonder why, but to do and analyze.)

Figure 4-8:
Residual
plots for
textbook-
weight data
minus the
outlier.



Knowing the Limitations of Your Regression Analysis

The bottom line of any data analysis is to make the correct conclusions given your results. When you're working with a simple linear regression model, there's the potential to make three major errors. This section shows you those errors and tells you how to

avoid them.

Avoiding slipping into cause-and-effect mode

In a simple linear regression, you investigate whether x is related to y , and if you get a strong correlation and a scatterplot that shows a linear trend, then you find the best-fitting line and use it to estimate the value of y for reasonable values of x .



There's a fine line, however (no pun intended), that you don't want to cross with your interpretation of regression results. Be careful to not automatically interpret slope in a cause-and-effect mode when you're using the regression line to estimate the value of y using x . Doing so can result in a leap of faith that can send you into the frying pan. Unless you have used a controlled experiment to get the data, you can only assume that the variables are correlated; you can't really give a stone-cold guarantee about why they're related.

In the textbook-weight example, you estimate the average weight of the students' textbooks by using the students' average weight, but that doesn't mean increasing a particular child's weight causes his textbook weight to increase. For example, because of the strong positive correlation, you do know that students with lower weights are associated with lower total textbook weights, and students with higher weights tend to have higher textbook weights. But you can't take one particular third-grade student, increase his weight, and presto — suddenly his textbooks weigh more.

The variable underlying the relationship between a child's weight and the weight of his backpack is the grade level of the student from an academic standpoint; as grade level increases, so might the size and number of his books, as well as the homework coming home. Student grade level drives both student weight and textbook weight. In this situation, student grade level is what statisticians call a *lurking variable*; it's a variable that wasn't included in the model but is related to both the outcome and the response. A lurking variable confuses the issue of what's causing what to happen.



If the collected data was the result of a well-designed experiment that controls for possible confounding variables, you can establish a cause-and-effect relationship between x and y if they're strongly correlated. Otherwise, you can't establish such a relationship. (See your Stats I text or Statistics For Dummies for info regarding experiments.)

Extrapolation: The ultimate no-no

Plugging values of x into the model that fall outside of the reasonable boundaries of x is

called *extrapolation*. And one of my colleagues sums up this idea very well: “Friends don’t let friends extrapolate.”

When you determine a best-fitting line for your data, you come up with an equation that allows you to plug in a value for x and get a predicted value for y . In algebra, if you find the equation of a line and graph it, the line typically has an arrow on each end indicating it goes on forever in either direction. But that doesn’t work for statistical problems (because statistics represents the *real* world). When you’re dealing with real-world units like height, weight, IQ, GPA, house prices, and the weight of your statistics textbook, only certain numbers make sense.

So the first point is, don’t plug in values for x that don’t make any sense. For example, if you’re estimating the price of a house (y) by using its square footage (x), you wouldn’t think of plugging in a value of x like 10 square feet or 100 square feet, because houses simply aren’t that small.

You also wouldn’t think about plugging in values like 1,000,000 square feet for x (unless your “house” is the Ohio State football stadium or something). It wouldn’t make sense. Likewise, if you’re estimating tomorrow’s temperature using today’s temperature, negative numbers for x could possibly make sense, but if you’re estimating the amount of precipitation tomorrow given the amount of precipitation today, negative numbers for x (or y for that matter) don’t make sense.



Choose only reasonable values of x for which you try to make estimates about y — that is, look at the values of x for which your data was collected, and stay within those bounds when making predictions. In the textbook-weight example, the smallest average student weight is 48.5 pounds, and the largest average student weight is 142 pounds. Choosing student weights between 48.5 and 142 to plug in for x in the equation is okay, but choosing values less than 48.5 or more than 142 isn’t a good idea. You can’t guarantee that the same linear relationship (or any linear relationship for that matter) continues outside the given boundaries.

Think about it: If the relationship you found actually continued for any value of x , no matter how large, then a 250-pound lineman from OSU would have to carry $3.69 + 0.113 * 250 = 31.94$ pounds of books around in his backpack. Of course this would be easy for him, but what about the rest of us?

Sometimes you need more than one variable

A simple linear regression model is just what it says it is: simple. I don’t mean easy to work with, necessarily, but simple in the uncluttered sense. The model tries to estimate the value of y by only using one variable, x . However, the number of real-world situations that can be explained by using a simple, one-variable linear regression is small. Often one variable just can’t do all the predicting.

If one variable alone doesn't result in a model that fits well enough, you can try to add more variables. It may take many variables to make a good estimate for y , and you have to be careful in how you choose them. In the case of stock market prices, for example, they're still looking for that ultimate prediction model.

As another example, health insurance companies try to estimate how long you'll live by asking you a series of questions (each of which represents a variable in the regression model). You can't find one single variable that estimates how long you'll live; you must consider many factors: your health, your weight, whether or not you smoke, genetic factors, how much exercise you do each week, and the list goes on and on and on.

The point is that regression models don't always use just one variable, x , to estimate y . Some models use two, three, or even more variables to estimate y . Those models aren't called simple linear regression models; they're called *multiple linear regression models* because of their employment of multiple variables to make an estimate. (You explore multiple linear regression models in Chapter 5.)

Chapter 5

Multiple Regression with Two X Variables

In This Chapter

- ▶ Getting the basic ideas behind a multiple regression model
 - ▶ Finding, interpreting, and testing coefficients
 - ▶ Checking model fit
-

The idea of regression is to build a model that estimates or predicts one quantitative variable (y) by using at least one other quantitative variable (x). Simple linear regression uses exactly one x variable to estimate the y variable. (See Chapter 4 for all the information you need on simple linear regression.) *Multiple linear regression*, on the other hand, uses more than one x variable to estimate the value of y .

In this chapter, you see how multiple regression works and how to apply it to build a model for y . You see all the steps necessary for the process, including determining which x variables to include, estimating their contributions to the model, finding the best model, using the model for estimating y , and assessing the fit of the model. It may seem like a mountain of information, but you won't regress on the topic of regression if you take this chapter one step at a time.

Getting to Know the Multiple Regression Model

Before you jump right into using the multiple regression model, get a feel for what it's all about. In this section, you see the usefulness of multiple regression as well as the basic elements of the multiple regression model. Some of the ideas are just an extension of the simple linear regression model (see Chapter 4). Some of the concepts are a little more complex, as you may guess because the model is more complex. But the concepts and the results should make intuitive sense, which is always good news.

Discovering the uses of multiple regression

One situation in which multiple regression is useful is when the y variable is hard to track down — that is, its value can't be measured straight up, and you need more than one other piece of information to help get a handle on what its value will be. For example, you may want to estimate the price of gold today. It would be hard to imagine being able to do that with only one other variable. You may base your estimate on recent gold prices, the price of other commodities on the market that move with or against gold, and a host of other possible economic conditions associated with the price of gold.

Another case for using multiple regression is when you want to figure out what factors play a role in determining the value of y . For example, you want to find out what information is important to real estate agents in setting a price for a house going on the market.

Looking at the general form of the multiple regression model

The general idea of simple linear regression is to fit the best straight line through that data that you possibly can and use that line to make estimates for y based on certain x -values. The equation of the best-fitting line in simple linear regression is $y = b_0 + b_1x_1$, where b_0 is the y -intercept and b_1 is the slope. (The equation also has the form $y = a + bx$; see Chapter 4.)

In the multiple regression setting, you have more than one x variable that's related to y . Call these x variables x_1, x_2, \dots, x_k . In the most basic multiple regression model, you use some or all of these x variables to estimate y where each x variable is taken to the first power. This process is called finding the best-fitting linear function for the data. This linear function looks like the following: $y = b_0 + b_1x_1 + b_2x_2 + \dots + b_kx_k$, and you can call it the *multiple (linear) regression model*. You use this model to make estimates about y based on given values of the x variables.



A *linear* function is an equation whose x terms are taken to the first power only.

For example $y = 2x_1 + 3x_2 + 24x_3$ is a linear equation using three x variables. If any of the x terms are squared, the function is a *quadratic* one; if an x term is taken to the third power, the function is a *cubic* function, and so on. In this chapter, I consider only linear functions.

Stepping through the analysis

Your job in conducting a multiple regression analysis is to do the following (the computer can help you do steps three through six):

- 1. Come up with a list of possible x variables that may be helpful in estimating y .**
- 2. Collect data on the y variable and your x variables from step one.**
- 3. Check the relationships between each x variable and y (using scatterplots and correlations), and use the results to eliminate those x variables that aren't strongly related to y .**
- 4. Look at possible relationships between the x variables to make sure you aren't being redundant (in statistical terms, you're trying to avoid the problem of multicollinearity).**

If two x variables relate to y the same way, you don't need both in the model.

5. Use those x variables (from step four) in a multiple regression analysis to find the best-fitting model for your data.

6. Use the best-fitting model (from step five) to predict y for given x -values by plugging those x -values into the model.

I outline each of these steps in the sections to follow.

Looking at x 's and y 's

The first step of a multiple regression analysis comes way before the number crunching on the computer; it occurs even before the data is collected. Step one is where you sit down and think about what variables may be useful in predicting your response variable y . This step will likely take more time than any other step, except maybe the data-collection process. Deciding which x variables may be candidates for consideration in your model is a deal-breaking step, because you can't go back and collect more data after the analysis is over.



Always check to be sure that your response variable, y , and at least one of the x variables are quantitative. For example, if y isn't quantitative but at least one x is, a logistic regression model may be in order (see Chapter 8).

Suppose you're in the marketing department for a major national company that sells plasma TVs. You want to sell as many TVs as you can, so you want to figure out which factors play a role in plasma TV sales. In talking with your advertising people and remembering what you learned in those business classes in college, you know that one powerful way to get sales is through advertising. You think of the types of advertising that may be related to sales of plasma TVs and your team comes up with two ideas:

- ✓ **TV ads:** Of course, how better to sell a TV than through a TV ad?
- ✓ **Newspaper sales:** Hit 'em on Sunday when they're reading the paper before watching the game through squinty eyes that are missing all the good plays and the terrible calls the referees are making.

By coming up with a list of possible x variables to predict y , you have just completed step one of a multiple regression analysis, according to the list in the previous section. Note that all three variables I use in the TV example are quantitative (the TV ad and newspaper sales variables and the TV sales response variable), which means you can go ahead and think about a multiple regression model by using the two types of ads to predict TV sales.

Collecting the Data

Step two in the multiple regression analysis process is to collect the data for your x and y variables. To do this, make sure that for each individual in the data set, you collect all the data for that individual at the same time (including the y -value and all x -values) and keep the data all together for each individual, preserving any relationships that may exist between the variables. You must then enter the data into a table format by using Minitab or any other software package (each column represents a variable and each row represents all the data from a single individual) to get a glimpse of the data and to organize it for later analyses.

To continue with the TV sales example from the preceding section, suppose that you start thinking about all the reams of data you have available to you regarding the plasma TV industry. You remember working with the advertising department before to do a media blitz by using, among other things, TV and newspaper ads. So you have data on these variables from a variety of store locations. You take a sample of 22 store locations in different parts of the country and put together the data on how much money was spent on each type of advertising, along with the plasma TV sales for that location. You can see the data in Table 5-1.

Table 5-1 Advertising Dollars and Sales of Plasma TVs

<i>Location</i>	<i>Sales (In Millions of Dollars)</i>	<i>TV Ads (In Thousands of Dollars)</i>	<i>Newspaper Ads (In Thousands of Dollars)</i>
1	9.73	0	20
2	11.19	0	20
3	8.75	5	5
4	6.25	5	5
5	9.10	10	10
6	9.71	10	10
7	9.31	15	15
8	11.77	15	15
9	8.82	20	5
10	9.82	20	5
11	16.28	25	25
12	15.77	25	25
13	10.44	30	0
14	9.14	30	0
15	13.29	35	5
16	13.30	35	5
17	14.05	40	10
18	14.36	40	10
19	15.21	45	15
20	17.41	45	15
21	18.66	50	20
22	17.17	50	20

In reviewing this data, the question is whether the amount of money spent on these two forms of advertising can do a good job of estimating sales (in other words, are the ads worth the money?). And if so, do you need to include spending for both types of ads to

estimate sales, or is one of them enough? Looking at the numbers in Table 5-1, you can see that higher sales may be related at least to higher amounts spent on TV advertising; the situation with newspaper advertising may not be so clear. So will the final multiple regression model contain both x variables or only one? In the following sections, you can find out.

Pinpointing Possible Relationships

The third step in doing a multiple regression analysis (see the list in the “Stepping through the analysis” section) is to find out which (if any) of your possible x variables are actually related to y . If an x variable has no relationship with y , including it in the model is pointless. Data analysts use a combination of scatterplots and correlations to examine relationships between pairs of variables (as you can see in Chapter 4). Although you can view these two techniques under the heading of looking for relationships, I walk you through each one separately in the following sections to discuss their nuances.

Making scatterplots

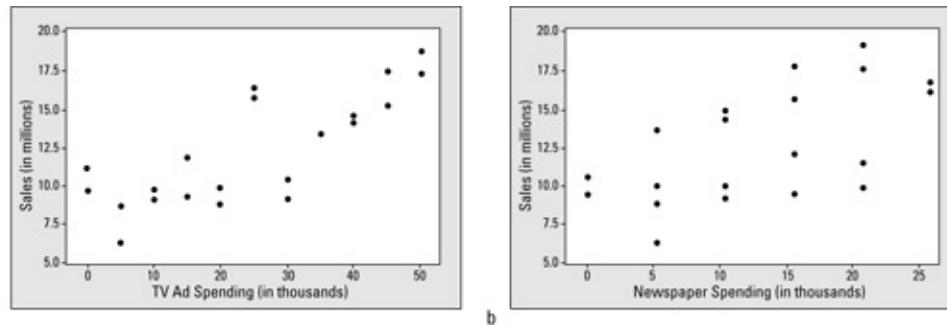
You make scatterplots in multiple linear regression to get a handle on whether your possible x variables are even related to the y variable you’re studying. To investigate these possible relationships, you make one scatterplot of each x variable with the response variable y . If you have k different x variables being considered for the final model, you make k different scatterplots.



To make a scatterplot in Minitab, enter your data in columns, where each column represents a variable and each row represents all the data from one individual. Go to Graph>Scatterplots>Simple. Select your y variable on the left-hand side, and click Select. That variable appears in the y -variable box on the right-hand side. Then select your x variable on the left-hand side, and click Select. That variable appears in the x -variable box on the right-hand side. Click OK.

Scatterplots of TV ad spending versus TV sales and newspaper ad spending versus TV sales are shown in Figure 5-1.

Figure 5-1:
Scatterplots
of TV and
newspaper
ad spending
versus
plasma TV
sales.



You can see from Figure 5-1a that TV spending does appear to have a fairly strong linear

relationship with sales. This observation provides evidence that TV ad spending may be useful in estimating plasma TV sales. Figure 5-1b shows a linear relationship between newspaper ad spending and sales, but the relationship isn't as strong as the one between TV ads and sales. However, it still may be somewhat helpful in estimating sales.

Correlations: Examining the bond

The second portion of step three involves calculating and examining the correlations between the x variables and the y variable. (Of course, if a scatterplot of an x variable and the y variable fails to come up with a pattern, then you drop that x variable altogether and don't proceed to find the correlation.)

Whenever you employ scatterplots to explore possible linear relationships, correlations are typically not far behind. The *correlation coefficient* is a number that measures the strength and direction of the linear relationship between two variables, x and y . (See Chapter 4 for the lowdown on correlation.)

This step involves two parts:

- ✓ Finding and interpreting the correlations
- ✓ Testing the correlations to see which ones are statistically significant (thereby determining which x variables are significantly related to y)

Finding and interpreting correlations

You can calculate a set of all possible correlations between all pairs of variables — which is called a *correlation matrix* — in Minitab. You can see the correlation matrix output for the TV data from Table 5-1 in Figure 5-2. Note the correlations between the y variable (sales) and each x variable, as well as the correlation between TV ads and newspaper ads.

Figure 5-2:
Correlation
values and p -
values for the
TV sales
example.

Correlations: Sales, TV, Newspaper			
TV	Sales	TV	
		0.791	
		0.000	
Newspaper	Sales	0.594	0.058
		0.004	0.799



Minitab can find a correlation matrix between any pairs of variables in the model, including the y variable and all the x variables as well. To calculate a correlation matrix for a group of variables in Minitab, first enter your data in columns (one for each variable). Then go to Stat>Basic Statistics>Descriptive Statistics>Correlation. Highlight the variables from the left-hand side for which you want correlations, and click Select.



To find the values of the correlation matrix from the computer output, intersect the row and column variables for which you want to find the correlation, and the top number in that intersection is the correlation of those two variables. For example, the correlation between TV ads and TV sales is 0.791, because it intersects the TV row with the Sales column in the correlation matrix in Figure 5-2.

Testing correlations for significance

By the rule-of-thumb approach from Stats I (also reviewed in Chapter 4), a correlation that's close to 1 or -1 (starting around ± 0.75) is strong; a correlation close to 0 is very weak/nonexistent; and around ± 0.6 to 0.7, the relationships become moderately strong. The correlation between TV ads and TV sales of 0.791 indicates a fairly strong positive linear relationship between these two variables, based on the rule-of-thumb. The correlation between newspaper ads and TV sales seen in Figure 5-2 is 0.594, which is moderate by my rule-of-thumb.

Many times in statistics a rule-of-thumb approach to interpreting a correlation coefficient is sufficient. However, you're in the big leagues now, so you need a more precise tool for determining whether or not a correlation coefficient is large enough to be statistically significant. That's the real test of any statistic: not that the relationship is fairly strong or moderately strong in the sample, but whether or not the relationship can be generalized to the population.

Now, that phrase *statistically significant* should ring a bell. It's your old friend the hypothesis test calling to you (see Chapter 3 for a brush-up on hypothesis testing). Just like a hypothesis test for the mean of a population or the difference in the means of two populations, you also have a test for the correlation between two variables within a population.



The null hypothesis to test a correlation is $H_0: \rho = 0$ (no relationship) versus $H_a: \rho \neq 0$ (a relationship exists). The letter ρ is the Greek version of r and represents the true correlation of x and y in the entire population; r is the correlation coefficient of the sample.

- ✓ **If you can't reject H_0 based on your data,** you can't conclude that the correlation between x and y differs from zero, indicating you don't have evidence that the two variables are related and x shouldn't be in the multiple regression model.
- ✓ **If you can reject H_0 based on your data,** you conclude that the correlation isn't equal to zero, so the variables are related. More than that, their relationship is deemed to be statistically significant — that is, the relationship would occur very rarely in your sample just by chance.

Any statistical software package can calculate a hypothesis test of a correlation for you. The actual formulas used in that process are beyond the scope of this book. However the interpretation is the same as for any test: If the p -value is smaller than your predetermined value of α (typically 0.05), reject H_0 and conclude x and y are related. Otherwise you can't reject H_0 , and you conclude you don't have enough evidence to indicate that the variables are related.



In Minitab, you can conduct a hypothesis test for a correlation by clicking on Stat>Basic Statistics>Correlation, and checking the Display p -values box. Choose the variables you want to find correlations for, and click Select. You'll get output in the form of a little table that shows the correlations between the variables for each pair with the respective p -values under each one. You can see the correlation output for the ads and sales example in Figure 5-2.

Looking at Figure 5-2, the correlation of 0.791 between TV ads and sales has a p -value of 0.000, which means it's actually less than 0.001. That's a highly significant result, much less than 0.05 (your predetermined α level). So TV ad spending is strongly related to sales. The correlation between newspaper ad spending and sales was 0.594, which is also found to be statistically significant with a p -value of 0.004.

Checking for Multicollinearity

You have one more very important step to complete in the relationship-exploration process before going on to using the multiple regression model. You need to complete step four: looking at the relationship between the x variables themselves and checking for redundancy. Failure to do so can lead to problems during the model-fitting process.



Multicollinearity is a term you use if two x variables are highly correlated. Not only is it redundant to include both related variables in the multiple regression model, but it's also problematic. The bottom line is this: If two x variables are significantly correlated, only include one of them in the regression model, not both. If you include both, the computer won't know what numbers to give as coefficients for each of the two variables because they share their contribution to determining the value of y . Multicollinearity can really mess up the model-fitting process and give answers that are inconsistent and often not repeatable in subsequent studies.

To head off the problem of multicollinearity, along with the correlations you examine regarding each x variable and the response variable y , also find the correlations between all pairs of x variables. If two x variables are highly correlated, don't leave them both in the model, or multicollinearity will result. To see the correlations between all the x variables, have Minitab calculate a correlation matrix of all the variables (see the section

“Finding and interpreting correlations”). You can ignore the correlations between the y variable and the x variables and only choose the correlations between the x variables shown in the correlation matrix. Find those correlations by intersecting the rows and columns of the x variables for which you want correlations.



If two x variables x_1 and x_2 are strongly correlated (that is, their correlation is beyond +0.7 or -0.7), then one of them would do just about as good a job of estimating y as the other, so you don’t need to include them both in the model. If x_1 and x_2 aren’t strongly correlated, then both of them working together would do a better job of estimating sales than either variable alone.

For the ad-spending example, you have to examine the correlation between the two x variables, TV ad spending and newspaper ad spending, to be sure no multicollinearity is present. The correlation between these two variables (as you can see in Figure 5-2) is only 0.058. You don’t even need a hypothesis test to tell you whether or not these two variables are related; they’re clearly not.

The p -value for the correlation between the spending for the two ad types is 0.799 (see Figure 5-2), which is much, much larger than 0.05 ever thought of being and therefore isn’t statistically significant. The large p -value for the correlation between spending for the two ad types confirms your thoughts that both variables together may be helpful in estimating y because each makes its own contribution. It also tells you that keeping them both in the model won’t create any multicollinearity problems. (This completes step four of the multiple regression analysis, as listed in the “Stepping through the analysis” section.)

Finding the Best-Fitting Model for Two x Variables

After you have a group of x variables that are all related to y and not related to each other (refer to previous sections), you’re ready to perform step five of the multiple regression analysis (as listed in the “Stepping through the analysis” section). You’re ready to find the best-fitting model for the data.

In the multiple regression model with two x variables, you have the general equation $y = b_0 + b_1x_1 + b_2x_2$, and you already know which x variables to include in the model (by doing step four in the previous section); the task now is to figure out which coefficients (numbers) to put in for b_0 , b_1 , and b_2 , so you can use the resulting equation to estimate y . This specific model is the *best-fitting multiple linear regression model*. This section tells you how to get, interpret, and test those coefficients in order to complete step five in the multiple regression analysis.



Finding the best-fitting linear equation is like finding the best-fitting line in simple linear regression, except that you're not finding a line. When you have two x variables in multiple regression, for example, you're estimating a best-fitting plane for the data.

Getting the multiple regression coefficients

In the simple linear regression model, you have the straight line $y = b_0 + b_1x$; the coefficient of x is the slope, and it represents the change in y per unit change in x . In a multiple linear regression model, the coefficients b_1 , b_2 , and so on quantify in a similar matter the sole contribution that each corresponding x variable (x_1 , x_2) makes in predicting y . The coefficient b_0 indicates the amount by which to adjust all these values in order to provide a final fit to the data (like the y -intercept does in simple linear regression).

Computer software does all the nitty-gritty work for you to find the proper coefficients (b_0 , b_1 , and so on) that fit the data best. The coefficients that Minitab settles on to create the best-fitting model are the ones that, as a group, minimize the sum of the squared residuals (sort of like the variance in the data around the selected model). The equations for finding these coefficients by hand are too unwieldy to include in this book; a computer can do all the work for you. The results appear in the regression output in Minitab. You can find the multiple regression coefficients (b_0 , b_1 , b_2 , \dots , b_k) on the computer output under the column labeled *COEF*.



To run a multiple regression analysis in Minitab, click on

Stat>Regression>Regression. Then choose the response variable (y) and click on Select. Then choose your predictor variables (x variables), and click Select. Click on OK, and the computer will carry out the analysis.

For the plasma TV sales example from the previous sections, Figure 5-3 shows the multiple regression coefficients in the *COEF* column for the multiple regression model. The first coefficient (5.257) is just the constant term (or b_0 term) in the model and isn't affiliated with any x variable. This constant just sort of goes along for the ride in the analysis; it's the number that you tack on the end to make the numbers work out right. The second coefficient in the *COEF* column is 0.162; this value is the coefficient of the x_1 (TV ad amount) term, also known as b_1 . The third coefficient in the *COEF* column is 0.249, which is the value for b_2 in the multiple regression model and is the coefficient that goes with x_2 (newspaper ad amount).

The regression equation is					
Sales = 5.267 + 0.162 TV ads + 0.249 Newsp ads					
Predictor	Coef	SE Coef	T	P	
Constant	5.2574	0.4984	10.55	0.000	
TV ads	0.16211	0.01319	12.29	0.000	
Newsp ads	0.24887	0.02792	8.91	0.000	
S = 0.976613	R-Sq = 92.8%	R-Sq(adj) = 92.0%			

Putting these coefficients into the multiple regression equation, you see the regression equation is Sales = 5.267 + 0.162 (TV ads) + 0.249 (Newspaper ads), where sales are in millions of dollars and ad spending is in thousands of dollars.

So you have your coefficients (no sweat, right?), but where do you go from here? What does it all mean? The next section guides you through interpretation.

Interpreting the coefficients

In simple linear regression (covered in Chapter 4), the coefficients represent the slope and y -intercept of the best-fitting line and are straightforward to interpret. The slope in particular represents the change in y due to a one-unit increase in x because you can write any slope as a number over one (and *slope* is rise over run).

In the multiple regression model, the interpretation's a little more complicated. Due to all the mathematical underpinnings of the model and how it's finalized (believe me, you don't want to go there unless you're looking for a PhD in statistics), the coefficients have a different meaning.



The coefficient of an x variable in a multiple regression model is the amount by which y changes if that x variable increases by one unit and the values of all other x variables in the model *don't change*. So basically, you're looking at the marginal contribution of each x variable when you hold the other variables in the model constant.

In the ads and sales regression analysis (see Figure 5-3), the coefficient of x_1 (TV ad spending) equals 0.16211. So y (plasma TV sales) increases by 0.16211 million dollars when TV ad spending increases by 1.0 thousand dollars and spending on newspaper ads doesn't change. (Note that keeping more digits after the decimal point reduces rounding error when in units of millions.)



You can more easily interpret the number "0.16211 million dollars" by converting it to a dollar amount without the decimal point: \$0.16211 million is equal to \$162,110. (To get this value, I just multiplied \$0.16211 by 1,000,000.) So plasma TV sales

increase by \$162,110 for each \$1,000 increase in TV ad spending and newspaper ad spending remains the same. Similarly, the coefficient of x_2 (newspaper ad spending) equals 0.24887. So plasma TV sales increase by 0.24887 million dollars (or \$248,870) when newspaper ad spending increases by \$1,000 and TV ad spending remains the same.



Don't forget the units of each variable in a multiple regression analysis. This mistake is one of the most common in Stats II. If you were to forget about units in the ads and sales example, you would think that sales increased by 0.24887 dollars with \$1 in newspaper ad spending!

Knowing the multiple regression coefficients (b_1 and b_2 , in this case) and their interpretation, you can now answer the original question: Is the money spent on TV or newspaper ads worth it? The answer is a resounding *yes!* Not only that, but you also can say how much you expect sales to increase per \$1,000 you spend on TV or newspaper advertising. Note that this conclusion assumes the model fits the data well. You have some evidence of that through the scatterplots and correlation tests, but more checking needs to be done before you can run to your manager and tell her the good news. The next section tells you what to do next.

Testing the coefficients

To officially determine whether you have the right x variables in your multiple regression model, do a formal hypothesis test to make sure the coefficients aren't equal to zero. Note that if the coefficient of an x variable is zero, when you put that coefficient into the model, you get zero times that x variable, which equals zero. This result is essentially saying that if an x variable's coefficient is equal to zero, you don't need that x variable in the model.



With any regression analysis, the computer automatically performs all the necessary hypothesis tests for the regression coefficients. Along with the regression coefficients you can find on the computer output, you see the test statistics and p -values for a test of each of those coefficients in the same row for each coefficient. Each one is testing H_0 : Coefficient = 0 versus H_a : Coefficient $\neq 0$.



The general format for finding a test statistic in most any situation is to take the statistic (in this case, the coefficient), subtract the value in H_0 (zero), and divide by the standard error of that statistic (for this example, the standard error of the coefficient). (For more info on the general format of hypothesis tests, see Chapter 3.)

To test a regression coefficient, the test statistic (using the labels from Figure 5-3) is $(\text{Coef} - 0)/\text{SE Coef}$. In noncomputer language, that means you take the coefficient, subtract zero, and divide by the standard error (SE) of the coefficient. The standard error of a coefficient here is a measure of how much the coefficient is expected to vary when you take a new sample. (Refer to Chapter 3 for more on standard error.)

The test statistic has a t -distribution with $n - k - 1$ degrees of freedom, where n equals the sample size and k is the number of predictors (x variables) in the model. This number of degrees of freedom works for any coefficient in the model (except you don't bother with a test for the constant, because it has no x variable associated with it).

The test statistic for testing each coefficient is listed in the column marked T (because it has a t -distribution) on the Minitab output. You compare the value of the test statistic to the t -distribution with $n - k - 1$ degrees of freedom (using Table A-1 in the appendix) and come up with your p -value. If the p -value is less than your predetermined α (usually 0.05), then you reject H_0 and conclude that the coefficient of that x variable isn't zero and that variable makes a significant contribution toward estimating y (given the other variables are also included in the model). If the p -value is larger than 0.05, you can't reject H_0 , so that x variable makes no significant contribution toward estimating y (when the other variables are included in the model).

In the case of the ads and plasma TV sales example, Figure 5-3 shows that the coefficient for the TV ads is 0.1621 (the second number in column two). The standard error is listed as being 0.0132 (the second number in column three). To find the test statistic for TV ads, take 0.1621 minus zero and divide by the standard error, 0.0132. You get a value of $t = 12.29$, which is the second number in column four. Comparing this value of t to a t -distribution with $n - k - 1 = 22 - 2 - 1 = 19$ degrees of freedom (Table A-1 in the appendix), you see the value of t is way off the scale. That means the p -value is smaller than can be measured on the t -table. Minitab lists the p -value in column five of Figure 5-3 as 0.000 (meaning it's less than 0.001). This result leads you to conclude that the coefficient for TV ads is statistically significant, and TV ads should be included in the model for predicting TV sales.

The newspaper ads coefficient is also significant with a p -value of 0.000 by the same reasoning; you find these results by looking across the newspaper ads row of Figure 5-3. Based on your coefficient tests and the lack of multicollinearity between TV and newspaper ads (see the earlier section “Checking for Multicollinearity”), you should include both the TV ads variable and the newspaper ads variable in the model for estimating TV sales.

Predicting y by Using the x Variables

When you have your multiple regression model, you're finally ready to complete step six of the multiple regression analysis: to predict the value of y given a set of values for the x variables. To make this prediction, you take those x values for which you want to predict

y , plug them into the multiple regression model, and simplify.

In the ads and plasma TV sales example (see analysis from Figure 5-3), the best-fitting model is $y = 5.26 + 0.162x_1 + 0.249x_2$. In the context of the problem, the model is Sales = $5.26 + 0.162$ TV ad spending (x_1) + 0.249 newspaper ad spending (x_2).



Remember that the units for plasma TV sales is in millions of dollars and the units for ad spending for both TV and newspaper ads is in the thousands of dollars. That is, \$20,000 spent on TV ads means $x_1 = 20$ in the model. Similarly, \$10,000 spent on newspaper ads means $x_2 = 10$ in the model. Forgetting the units involved can lead to serious miscalculations.

Suppose you want to estimate plasma TV sales if you spend \$20,000 on TV ads and \$10,000 on newspaper ads. Plug $x_1 = 20$ and $x_2 = 10$ into the multiple regression model, and you get $y = 5.26 + 0.162(20) + 0.249(10) = 10.99$. In other words, if you spend \$20,000 on TV advertising and \$10,000 in newspaper advertising, you estimate that sales will be \$10.99 million.

This estimate at least makes sense in terms of the data from the 22 store locations shown in Table 5-1. Location 10 spent \$20,000 on TV ads and \$5,000 on newspaper ads (short of what you had) and got sales of \$9.82 million. Location 11 spent a little more on TV ads and a lot more on newspaper ads than what you had and got sales of \$16.28 million. Your estimates of sales for Store Locations 10 and 11 are $5.26 + 0.162 * 20 + 0.249 * 5 = \9.745 million, and $5.26 + 0.162 * 25 + 0.249 * 25 = \15.535 million, respectively. These estimates turned out to be pretty close to the actual sales at those two locations (\$9.82 million and \$16.28 million, respectively, as shown in Table 5-1), giving at least some confidence that your estimates will be close for the other store locations not chosen for the study.



Be careful to put in only values for the x variables that fall in the range of where the data lies. In other words, Table 5-1 shows data for TV ad spending between \$0 and \$50,000; newspaper ad spending goes from \$0 to \$25,000. It wouldn't be appropriate to try to estimate sales for spending amounts of \$75,000 for TV ads and \$50,000 for newspaper ads, respectively, because the regression model you came up with only fits the data that you collected. You have no way of knowing whether that same relationship continues outside that area. This no-no of estimating y for values of the x variables outside their range is called *extrapolation*. As one of my colleagues says, "Friends don't let friends extrapolate."

Checking the Fit of the Multiple Regression Model

Before you run to your boss in triumph saying you've slam-dunked the question of how to estimate plasma TV sales, you first have to make sure all your i's are dotted and all your t's are crossed, as you do with any other statistical procedure. In this case, you have to check the conditions of the multiple regression model. These conditions mainly focus on the *residuals* (the difference between the estimated values for y and the observed values of y from your data). If the model is close to the actual data you collected, you can feel somewhat confident that if you were to collect more data, it would fall in line with the model as well, and your predictions should be good.

In this section, you see what the conditions are for multiple regression and specific techniques statisticians use to check each of those conditions. The main character in all this condition-checking is the residual.

Noting the conditions

The conditions for multiple regression concentrate on the error terms, or residuals. The residuals are the amount that's left over after the model has been fit. They represent the difference between the actual value of y observed in the data set and the estimated value of y based on the model. Following are the conditions for the residuals of the multiple regression model; note that all conditions need to be met in order to give the go-ahead for a multiple regression model:

- ✓ They have a normal distribution with a mean of zero.
- ✓ They have the same variance for each fitted (predicted) value of y .
- ✓ They're independent (meaning they don't affect each other).

Plotting a plan to check the conditions

It may sound like you have a ton of things to check here and there, but luckily, Minitab gives you all the info you need to know in a series of four graphs, all presented at one time. These plots are called the *residual plots*, and they graph the residuals so that you can check to see whether the conditions from the previous section are met.

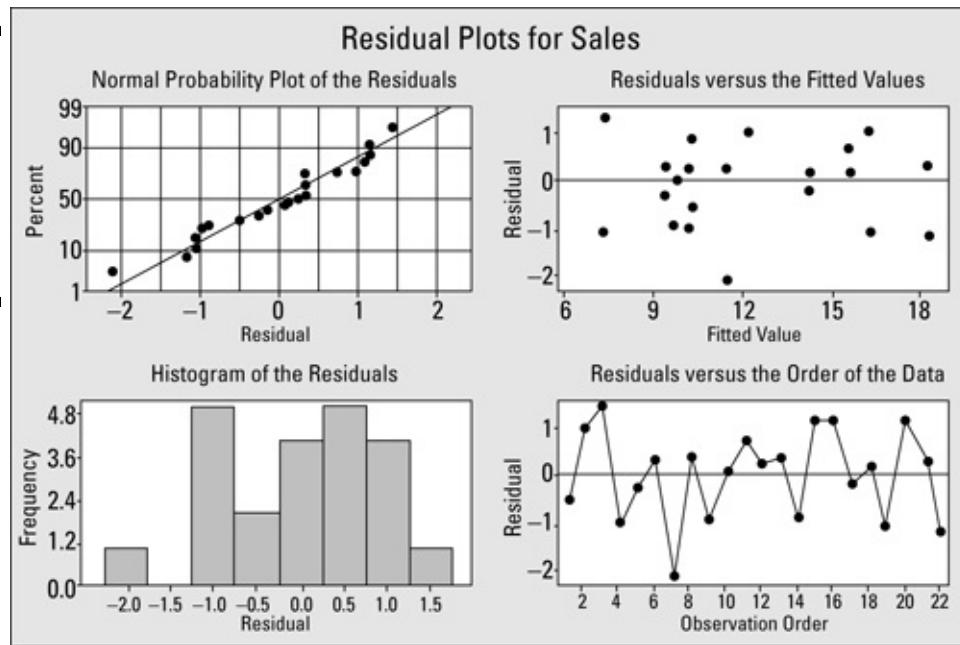
You can get the set of residual plots in two flavors:

- ✓ **Regular residuals:** The regular residual plots (the vanilla-flavored ones) show you exactly what the residuals are for each value of y . Their units depend on the variables in the model; use them only if you want to mainly look for patterns in the data. Figure 5-4 shows the plots of the regular residuals for the TV sales example. These residuals are in units of millions of dollars.
- ✓ **Standardized residuals:** The standardized residual plots (the strawberry-flavored kind) take each residual and convert it to a Z-score by subtracting the mean and dividing by the standard deviation of all the residuals. Figure 5-5 shows the plots of the standardized residuals for the TV sales example. Use these plots if you want to

not only look for patterns in the data but also assess the standardized values of the residuals in terms of values on a Z-distribution to check for outliers. (Most statisticians use standardized residual plots.)

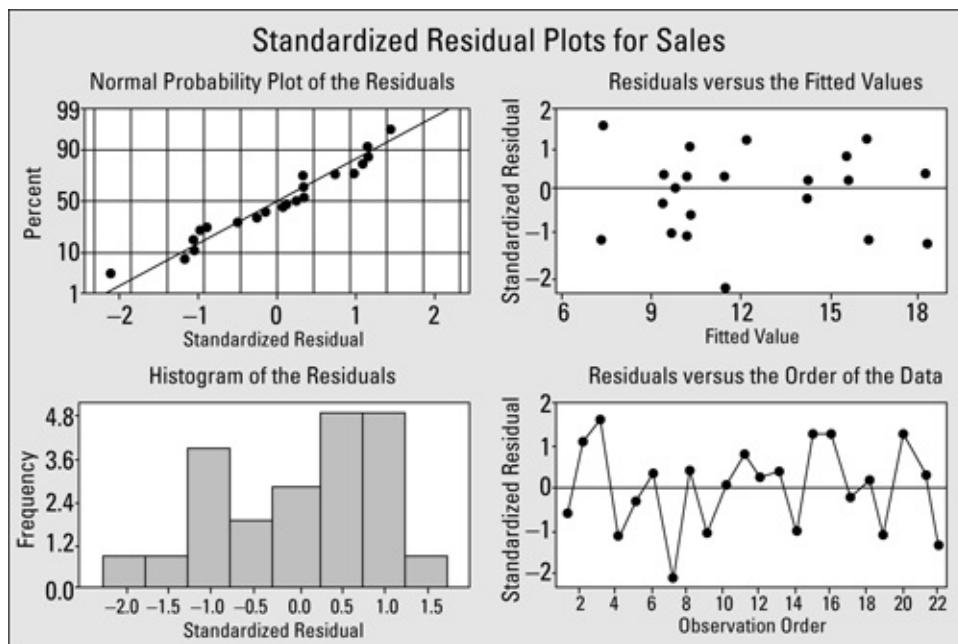
Note that the plots in Figure 5-5 look almost exactly the same as those in Figure 5-4. It's not surprising that the shapes of all graphs are the same for both types of residuals. Note however that the values of the regular residuals in Figure 5-4 are in millions of dollars and the standardized residuals in Figure 5-5 are from the standard normal distribution, which has no units.

Figure 5-4:
Residual plots for the ads and plasma TV sales example.



To make residual plots in Minitab, go to Stat>Regression>Regression. Select your response (y) variable and your predictor (x) variables. Click on Graphs, and choose either Regular or Standardized for the residuals, depending on which one you want. Then click on Four-in-one, which indicates you want to get all four residual plots shown in Figure 5-4 (using regular residuals) and Figure 5-5 (using standardized residuals).

Figure 5-5:
Standardized residual plots for the ads and plasma TV example.



Checking the three conditions

The following sections show you how to check the residuals to see whether your data set meets the three conditions of the multiple regression model.

Meeting the first condition: Normal distribution with mean zero

The first condition to meet is that the residuals must have a normal distribution with mean zero. The upper-left plot of Figure 5-4 shows how well the residuals match a normal distribution. Residuals falling in a straight line means the normality condition is met. By the looks of this plot, I'd say that condition is met for the ad and sales example.

The upper-right plot of Figure 5-4 shows what the residuals look like for the various estimated y values. Look at the horizontal line going across that plot: It's at zero as a marker. The residuals should average out to be at that line (zero). This Residuals versus Fitted Values plot checks the mean-of-zero condition and holds for the ads and sales example looking at Figure 5-4.



As an alternative check for normality apart from using the regular residuals, you can look at the standardized residuals plot (see Figure 5-5) and check out the upper-right plot. It shows how the residuals are distributed across the various estimated (fitted) values of y . Standardized residuals are supposed to follow a standard normal distribution — that is, they should have mean of zero and standard deviation of one. So when you look at the standardized residuals, they should be centered around zero in a way that has no predictable pattern, with the same amount of variability around the horizontal line that crosses at zero as you move from left to right.

In looking at the upper-right plot of Figure 5-5, you should also find that most (95

percent) of the standardized residuals fall within two standard deviations of the mean, which in this case is -2 to $+2$ (via the 68-95-99.7 Rule — remember that from Stats I?). You should see more residuals hovering around zero (where the middle lump would be on a standard normal distribution), and you should have fewer and fewer of the residuals as you go away from zero. The upper-right plot in Figure 5-5 confirms a normal distribution for the ads and sales example on all the counts mentioned here.

The lower-left plots of Figures 5-4 and 5-5 show histograms of the regular and standardized residuals, respectively. These histograms should reflect a normal distribution; the shape of the histograms should be approximately symmetric and look like a bell-shaped curve. If the data set is small (as is the case here with only 22 observations), the histogram may not be as close to normal as you would like; in that case, consider it part of the body of evidence that all four residual plots show you. The histograms shown in the lower-left plots of Figure 5-4 and 5-5 aren't terribly normal looking; however, because you can't see any glaring problems with the upper-right plots, don't be worried.

Satisfying the second condition: Variance

The second condition in checking the multiple regression model is that the residuals have the same variance for each fitted (predicted) value of y . Look again at the upper-right plot of Figure 5-4 (or Figure 5-5). You shouldn't see any change in the amount of spread (variability) in the residuals around that horizontal line as you move from left to right. Looking at the upper-right graph of Figure 5-4, there's no reason to say condition number two hasn't been met.



One particular problem that raises a red flag with the second condition is if the residuals fan out, or increase in spread, as you move from left to right on the upper-right plot. This fanning out means that the variability increases more and more for higher and higher predicted values of y , so the condition of equal variability around the fitted line isn't met, and the regression model wouldn't fit well in that case.

Checking the third condition

The third condition is that the residuals are independent; in other words, they don't affect each other. Looking at the lower-right plot on either Figure 5-4 or 5-5, you can see the residuals plotted by *observation number*, which is the order in which the data came in the sample. If you see a pattern, you have trouble; for example, if you were to connect the dots, so to speak, you might see a pattern of a straight line, a curve, or any kind of predictable up or down trend. You can see no patterns in the lower-right plots, so the independence condition is met for the ads and plasma TV sales example.



If the data must be collected over time, such as stock prices over a ten-year period, the independence condition may be a big problem because the data from the previous time period may be related to the data from the next time period. This kind of data requires time series analysis and is beyond the scope of this book.

Chapter 6

How Can I Miss You If You Won't Leave? Regression Model Selection

In This Chapter

- ▶ Evaluating different methods for choosing a multiple regression model
- ▶ Understanding how forward selection and backward selection work
- ▶ Using the best subsets methods to find a good model

Suppose you're trying to estimate some quantitative variable, y , and you have many x variables available at your disposal. You have so many variables related to y , in fact, that you feel like I do in my job every day — overwhelmed with opportunity. Where do you go? What do you do? Never fear, this chapter is for you.

In this chapter, you uncover criteria for determining when a model fits well. I discuss different model selection procedures and all the details of the most statistician-approved method for selecting the best model. Plus, you get to find out what factors come into play when a punter kicks a football. (You can think about that while you're reading.)



Note that the term *best* has many connotations here. You can't find one end-all-be-all model that everyone comes up with in the end. That's to say that each data analyst can come up with a different model, and each model still could do a good job of predicting y .

Getting a Kick out of Estimating Punt Distance

Before you jump into a model selection procedure to predict y by using a set of x variables, you have to do some legwork. The variable of interest is y , and that's a given. But where do the x variables come from? How do you choose which ones to investigate as being possible candidates for predicting y ? And how do those possible x variables interact with each other toward making that prediction?

You must answer all these questions before using any model selection procedure. However, this part is the most challenging and the most fun; a computer can't think up x variables for you!

Suppose you're at a football game and the opposing team has to punt the ball. You see the punter line up and get ready to kick the ball, and some questions come to you: "Gee, I wonder how far this punt will go? I wonder what factors influence the distance of a punt?"

Can I use those factors in a multiple regression model to try to estimate punt distance? Hmm, I think I'll consult my *Statistics II For Dummies* book on this and analyze some data during halftime. . . .”

Well, maybe that's pushing it, but it's still an interesting line of questioning for football players, golfers, soccer players, and even baseball players. Everyone's looking for more distance and a way to get it.

In the following sections, you can see how to identify and assess different x variables in terms of their potential contribution to predicting y .

Brainstorming variables and collecting data

Starting with a blank slate and trying to think of a set of x variables that may be related to y may sound like a daunting task, but in reality, it's probably not as bad as you think. Most researchers who are interested in predicting some variable y in the first place have some ideas about which variables may be related to it. After you come up with a set of logical possibilities for x , you collect data on those variables, as well as on y , to see what their actual relationship with y may be.

The Virginia Polytechnic Institute did a study to try to estimate the distance of a punt in football (something Ohio State fans aren't familiar with). Possible variables they thought may be related to the distance of a punt included the following:

- ✓ Hang time (time in the air, in seconds)
- ✓ Right leg strength (measured in pounds of force)
- ✓ Left leg strength (in pounds of force)
- ✓ Right leg flexibility (in degrees)
- ✓ Left leg flexibility (in degrees)
- ✓ Overall leg strength (in pounds)

The data collected on a sample of 13 punts (by right-footed punters) is shown in Table 6-1.

Table 6-1 Data Collected for Punt Distance Study

Distance (In Feet)	Hang Time	Right Leg Strength	Left Leg Strength	Right Leg Flexibility	Left Leg Flexibility	Overall Leg Strength
162.50	4.75	170	170	106	106	240.57
144.00	4.07	140	130	92	93	195.49
147.50	4.04	180	170	93	78	152.99
163.50	4.18	160	160	103	93	197.09
192.00	4.35	170	150	104	93	266.56
171.75	4.16	150	150	101	87	260.56
162.00	4.43	170	180	108	106	219.25
104.93	3.20	110	110	86	92	132.68
105.67	3.02	120	110	90	86	130.24
117.59	3.64	130	120	85	80	205.88
140.25	3.68	120	140	89	83	153.92
150.17	3.60	140	130	92	94	154.64
165.17	3.85	160	150	95	95	240.57

Other variables you may think of that are related to punt distance may include the direction and speed of the wind at the time of the punt, the angle at which the ball was snapped, the average distance of punts made in the past by a particular punter, whether the game is at home or away in a hostile environment, and so on. However, these researchers seem to have enough information on their hands to build a model to estimate punt distance.

For the sake of simplicity, you can assume the kicker is right-footed, which isn't always the case, but it represents the overwhelming majority of kickers.

Looking just at this raw data set in Table 6-1, you can't figure out which variables, if any, are related to distance of the punt or how those variables may be related to punt distance. You need more analyses to get a handle on this.

Examining scatterplots and correlations



After you've identified a set of possible x variables, the next step is to find out which of these variables are highly related to y in order to start trimming down the set of possible candidates for the final model. In the punt distance example, the goal is to see which of the six variables in Table 6-1 are strongly related to punt distance. The two ways to look at these relationships are

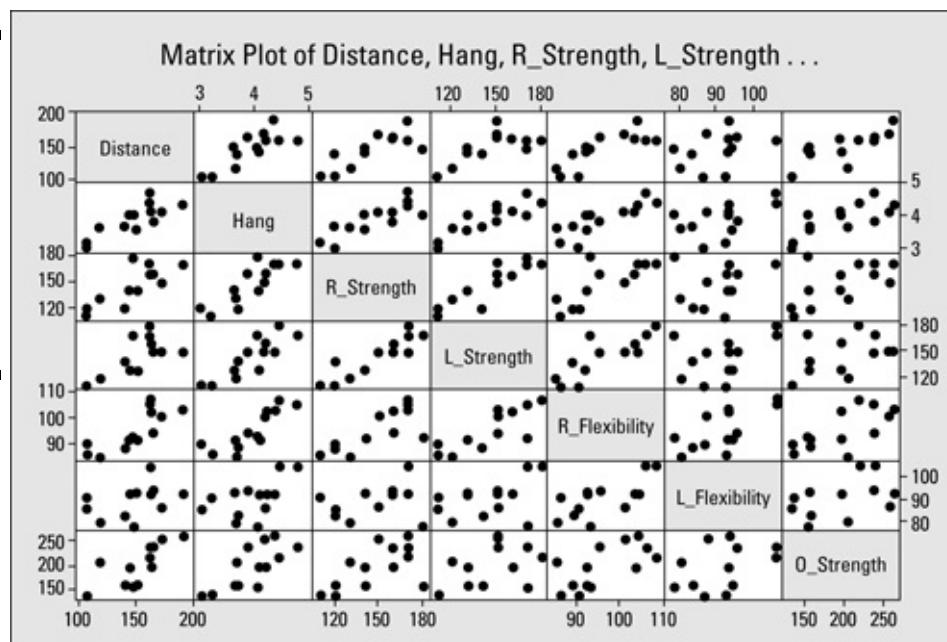
- ✓ **Scatterplot:** A graphical technique
- ✓ **Correlation:** A one-number measure of the linear relationship between two variables

Seeing relationships through scatterplots

To begin examining the relationships between the x variables and y , you use a series of scatterplots. Figure 6-1 shows all the scatterplots — not only of each x variable with y but also of each x variable with the other x variables. The scatterplots are in the form of a *matrix*, which is a table made of rows and columns. For example, the first scatterplot in row two of Figure 6-1 looks at the variables of distance (which appears in column one) and hang time (which appears in row two). This scatterplot shows a possible positive (uphill) linear relationship between distance and hang time.

Note that Figure 6-1 is essentially a symmetric matrix across the diagonal line. The scatterplot for distance and hang time is the same as the scatterplot for hang time and distance; the x and y axes are just switched. The essential relationship shows up either way. So you only have to look at all the scatterplots below the diagonal (where the variable names appear) or above the diagonal. You don't need to examine both.

Figure 6-1: A matrix of all scatter-plots between pairs of variables in the punting distance example.



To get a matrix of all scatterplots between a set of variables in Minitab, go to Graph>Matrix Plot and choose Matrix of Plots>Simple. Highlight all the variables in the left-hand box for which you want scatterplots by clicking on them; click Select, and then click OK. You'll see the matrix of scatterplots with a format similar to Figure 6-1.

Looking across row one of Figure 6-1, you can see that all the variables seem to have a positive linear relationship with punt distance except left leg flexibility. Perhaps the reason left leg flexibility isn't much related to punt distance is because the left foot is planted into the ground when the kick is made — for a right-footed kicker, the left leg doesn't have to be nearly as flexible as the right leg, which does the kicking. So it doesn't appear that left leg flexibility contributes a great deal to the estimation of punt distance on its own.

You can also see in Figure 6-1 that the scatterplots showing relationships between pairs of x variables are to the right of column one and below row one. (Remember, you need to look on only the bottom part of the matrix or the top part of the matrix to see the relevant scatterplots.) It appears that hang time is somewhat related to each of the other variables (except left leg flexibility, which doesn't contribute to estimating y). So hang time could possibly be the most important single variable in estimating the distance of a punt.

Looking for connections by using correlations

Scatterplots can give you some general ideas as to whether two variables are related in a linear way. However, pinpointing that relationship requires a numerical value to tell you how strongly the variables are related (in a linear fashion) as well as the direction of that relationship. That numerical value is the *correlation* (also known as *Pearson's correlation*; see Chapter 4). So the next step toward trimming down the possible candidates for x variables is to calculate the correlation between each x variable and y .



To get a set of all the correlations between any set of variables in your model by using Minitab, go to Stat>Basic Statistics>Correlation. Then highlight all the variables you want correlations for, and click Select. (To include the p -values for each correlation, click the Display p -values box.) Then click OK. You'll see a listing of all the variables' names across the top row and down the first column. Intersect the row depicting the first variable with the column depicting the second variable in order to find the correlation for that pair.

Table 6-2 shows the correlations you can calculate between y = punt distance and each of the x variables. These results confirm what the scatterplots were telling you. Distance seems to be related to all the variables except left leg flexibility because that's the only variable that didn't have a statistically significant correlation with distance using the α level 0.05. (For more on the test for correlation, see Chapter 5.)

Table 6-2 Correlations between Distance of a Punt and Other Variables

<i>x</i> Variable	Correlation with Punt Distance	<i>p</i> -value
Hang time	0.819	0.001*
Right leg strength	0.791	0.001*
Left leg strength	0.744	0.004*
Right leg flexibility	0.806	0.001*
Left leg flexibility	0.408	0.167
Overall leg strength	0.796	0.001*

* Statistically significant at level $\alpha = 0.05$

If you take a look at Figure 6-1, you can see that hang time is related to other x variables such as right foot and left foot strength, right leg flexibility, and so on. This is where

things start to get sticky. You have hang time related to distance, and lots of other variables related to hang time. Although hang time is clearly the most related to distance, the final multiple regression model may not include hang time.

Here's one possible scenario: You find a combination of other x variables that can do a good job estimating y together. And all those other variables are strongly related to hang time. This result may mean that in the end you don't need to include hang time in the model. Strange things happen when you have many different x variables to choose from.

After you narrow down the set of possible x variables for inclusion in the model to predict punt distance, the next step is to put those variables through a selection procedure to trim down the list to a set of essential variables for predicting y .

Just Like Buying Shoes: The Model Looks Nice, But Does It Fit?

When you get into model selection procedures, you find that many different methods exist for selecting the best model, according to a wide range of criteria. Each one can result in models that differ from each other, but that's something I love about statistics: Sometimes there's no one single best answer.

The three model selection procedures covered in this section are

- ✓ **Best subsets procedure**
- ✓ **Forward selection**
- ✓ **Backward selection**

Of all the model selection procedures out there, the one that gets the most votes with statisticians is the *best subsets procedure*, which examines every single possible model and determines which one fits best, using certain criteria.

In this section, you see different methods statisticians use to assess and compare the fit of different models. You see how the best subsets procedure works for model selection in a step-by-step manner. Then I show you how to take all the information given to you and wade through it to make your way to the answer — the best-fitting model based on a subset of the available x variables. Finally, you see how this procedure is applied to find a model to predict punt distance.

Assessing the fit of multiple regression models

For any model selection procedure, assessing the fit of each model being considered is built into the process. In other words, as you go through all the possible models, you're always keeping an eye on how well each model fits. So before you get into a discussion of

how to do the best subsets procedure, you need criteria to assess how well a particular model fits a data set.

Although there are tons of different statistics for assessing the fit of regression models, I discuss the most popular ones: R^2 (simple linear regression only), R^2 adjusted, and Mallow's C-p. All three appear on the bottom line of the Minitab output when you do any sort of model selection procedure. Here's a breakdown of the assessment techniques:

- ✓ **R^2 :** R^2 is the percentage of the variability in the y values that's explained by the model. It falls between 0 and 100 percent (0 and 1.0). In simple linear regression (see Chapter 4), a high value of R^2 means the line fits well, and a low value of R^2 means the line doesn't fit well.

When you have multiple regression, however, there's a bit of a catch here. As you add more and more variables (no matter how significant), the value of R^2 increases or stays the same — it never goes down. This can result in an inflated measure of how well the model fits. Of course, statisticians have a fix for the problem, which leads me to the next item on this list.

- ✓ **R^2 adjusted:** R^2 adjusted takes the value of R^2 and adjusts it downward according to the number of variables in the model. The higher the number of variables in the model, the lower the value of R^2 adjusted will be, compared to the original R^2 .

A high value of R^2 adjusted means the model you have is fitting the data very well (the closer to 1, the better). I typically find a value of 0.70 to be considered okay for R^2 adjusted, and the higher the better.



Always use R^2 adjusted rather than the regular R^2 to assess the fit of a multiple regression model. With every addition of a new variable into a multiple regression model, the value of R^2 stays the same or increases. It will never go down because a new variable will either help explain some of the variability in the y 's (thereby increasing R^2 by definition), or it will do nothing (leaving R^2 exactly where it was before). So theoretically, you could just keep adding more and more variables into the model just for the sake of getting a larger value of R^2 .

R^2 adjusted is important because it keeps you from adding more and more variables by taking into account how many variables there already are in the model. The value of R^2 adjusted can actually decrease if the added value of the additional variable is outweighed by the number of variables in the model. This gives you an idea of how much or how little added value you get from a bigger model (bigger isn't always better).

- ✓ **Mallow's C-p:** Mallow's C-p takes the amount of error left unexplained by a model of p with the x variables, divides that number by the average amount of error left over from the full model (with all the x variables), and adjusts that result for the number of observations (n) and the number of x variables used (p). In general, the smaller Mallow's C-p is, the better, because when it comes to the amount of error in your

model, less is more. A C-p value close to p (the number of x variables in the model) reflects a model that fits well.

Model selection procedures

The process of finding the “best” model is not cut and dry. (Heck, even the definition of “best” here isn’t cut and dry.) Many different procedures exist for going through different models in a systematic way, evaluating each one, and stopping at the right model. Three of the more common model selection procedures are forward selection, backward selection, and the best subsets model. In this section you get a very brief overview of the forward and backward selection procedures, and then you get into the details of the best subsets model, which is the one statisticians use most.

Going with the forward selection procedure

The forward selection procedure starts with a model with no variables in it and adds variables one at a time according to the amount of contribution they can make to the model.

Start with an entry level value of α . Then run hypothesis tests (see Chapter 3 for instructions) for each x variable to see how it’s related to y . The x variable with the smallest p-value wins and is added to the model, as long as its p-value is smaller than the entry level. You keep doing this with the remaining variables until the one with the smallest p-value doesn’t make the entry level. Then you stop.



The drawback of the forward selection procedure is that it starts with nothing and adds variables one at a time as you go along; after a variable is added, it’s never removed. The best model might not even get tested.

Opting for the backward selection procedure

The backward selection procedure does the opposite of the forward selection method. It starts with a model with all the x variables in it and removes variables one at a time. Those that make the least amount of contribution to the model are removed first. You choose a removal level to begin; then you test all the x variables and find the one with the largest p-value. If the p-value of this x variable is higher than the removal level, that variable is taken out of the model.

You continue removing variables from the model until the one with the largest p-value doesn’t exceed the removal level. Then you stop.



The drawback of the backward selection procedure is that it starts with

everything and removes variables one at a time as you go along; after a variable is removed, it never comes back. Again, the best model might not even be tested.

Using the best subsets procedure

The best subsets procedure has fewer steps than the forward or backward selection model because the computer formulates and analyzes all possible models in a single step. In this section, you see how to get the results and then use them to come up with a best multiple regression model for predicting y .

Here are the steps for conducting the best subsets model selection procedure to select a multiple regression model; note that Minitab does all the work for you to crunch the numbers:

1. Conduct the best subsets procedure in Minitab, using all possible subsets of the x variables being considered for inclusion in the final model.



To carry out the best subsets selection procedure in Minitab, go to Stat>Regression>Best Subsets. Highlight the response variable (y), and click Select. Highlight all the predictor (x) variables, click Select, and then click OK.

The output contains a listing of all models that contain one x variable, all models that contain two x variables, all models that contain three x variables, and so on, all the way up to the full model (containing all the x variables). Each model is presented in one row of the output.

2. Choose the best of all the models shown in the best subsets Minitab output by finding the model with the largest value of R^2 adjusted and the smallest value of Mallow's C-p; if two competing models are about equal, choose the model with the fewer number of variables.

If the model fits well, R^2 adjusted is high. So you also want to look for the smallest possible model that has a high value of R^2 adjusted and a small value of Mallow's C-p compared to its competitors. And if it comes down to two similar models, always make your final model as easy to interpret as possible by selecting the model with fewer variables.

The secret to a punter's success: An example

Returning to the punt distance example from earlier in this chapter, suppose that you analyzed the punt distance data by using the best subsets model selection procedure. Your results are shown in Figure 6-2. This section follows Minitab's footsteps in getting these results and provides you with a guide for interpreting the results.

Assuming that you already used Minitab to carry out the best subsets selection procedure on the punt distance data, you can now analyze the output from Figure 6-2. Each variable shows up as a column on the right side of the output. Each row represents

the results from a model containing the number of variables shown in column one. The X's at the end of each row tell you which variables were included in that model. The number of variables in the model starts at 1 and increases to 6 because six x variables are available in the data set.

The models with the same number of variables are ordered by their values of R^2 adjusted and Mallows C-p, from best to worst. The top-two models (for each number of variables) are included in the computer output.

For example, rows one and two of Figure 6-2 (both marked 1 in the Vars column) show the top-two models containing one x variable; rows three and four show the top two models containing two x variables; and so on. Finally, the last row shows the results of the full model containing all six variables. (Only one model contains all six variables, so you don't have a second-best model in this case.)

Looking at the first two rows of Figure 6-2, the top one-variable model is the one including hang time only. The second-best one-variable model includes only right foot flexibility. The right foot flexibility model has a lower value of R^2 and a higher Mallow's C-p than the hang time model, which is why it's the second best.

Row three shows that the best two-variable model for estimating punt distance is the model containing right leg strength and overall leg strength. The best three-variable model is in row five; it shows that the best three-variable model includes right foot strength, right foot flexibility, and overall leg strength. The best four-variable model is found in row seven and includes right foot strength, right and left foot flexibility, and overall foot strength. The best five-variable model is found in row nine and includes every variable except left foot strength. The only six-variable model with all variables included is listed in the last row.

Among the best one-variable, two-variable, three-variable, four-variable, and five-variable models, which one should you choose for your final multiple regression model? Which model is the best of the best? With all these results, it would be easy to have a major freakout over which one to pick, but never fear — Mallow's is here (along with his friendly sidekick, the R^2 adjusted).

Looking at Figure 6-2 column three, you see that as the number of variables in the model increases, R^2 adjusted peaks out and then drops way off. That's because R^2 adjusted takes into account the number of variables in the model and reduces R^2 accordingly. You can see that R^2 adjusted peaks out at a level of 74.1 percent for two models. The corresponding models are the top two-variable model (right leg strength and overall leg strength) and the best three-variable model (right foot strength, right foot flexibility, and overall leg strength).

Now look at Mallow's C-p for these two models. Notice that Mallow's C-p is zero for the best two-variable model and 1.3 for the best three-variable model. Both values are small compared to others in Figure 6-2, but because Mallow's C-p is smaller for the two-variable

model, and because it has one less variable in it, you should choose the two-variable model (right leg strength and overall leg strength) as the final model, using the best subsets procedure.

Figure 6-2:
Best subsets
procedure
results for the
punt distance
example.

Best Subsets Regression: Distance versus Hang, R_Strength ...						
Response is Distance						
Mallows	R	L				
			<u>F</u>	<u>F</u>		
	R	L	i	1	1	O
	S	S	x	x	S	
	t	t	i	i	t	
	r	r	b	b	r	
	e	e	i	i	e	
	H	n	n	l	l	n
	a	g	g	i	i	g
	n	t	t	t	t	t
Vars	R-Sq	R-Sq(adj)	C-p	S	g	h
	1	67.1	64.1	1.7	15.570	X
	1	65.0	61.8	2.3	16.043	X
	2	78.5	74.1	-0.0	13.206	X X
	2	78.2	73.8	0.1	13.294	X X
	3	80.6	74.1	1.3	13.214	X X X
	3	79.5	72.7	1.6	13.581	X X
	4	81.4	72.1	3.0	13.724	X X X X
	4	80.7	72.0	3.3	13.977	X X X X
	5	81.5	68.2	5.0	14.643	X X X X X
	5	81.4	68.2	5.0	14.650	X X X X X
	6	81.5	62.9	7.0	15.812	X X X X X X

Chapter 7

Getting Ahead of the Learning Curve with Nonlinear Regression

In This Chapter

- ▶ Getting a feel for nonlinear regression
 - ▶ Making use of scatterplots
 - ▶ Fitting a polynomial to your data set
 - ▶ Exploring exponential models to fit your data
-

In Stats I, you concentrate on the *simple linear regression model*, where you look for one quantitative variable, x , that you can use to make a good estimate of another quantitative variable, y , using a straight line. The examples you look at in Stats I fall right in line with this kind of model, such as using height to estimate weight or using study time to estimate exam score. (For more information and examples for using simple linear regression models, see Chapter 4.)

But not all situations fall into the straight line category. Take gas mileage and speed, for example. At low speeds, gas mileage is lower, and at high speeds, gas mileage is lower; but at medium speeds, gas mileage is higher. This low-high-low relationship between speed and gas mileage represents a curved relationship. Relationships that don't resemble straight lines are called *nonlinear relationships* (clever, huh?). Looked at simply, nonlinear regression takes the stage when you want to predict some quantitative variable (y) by using another quantitative variable (x) but the pattern you see in the data collected resembles a curve, not a straight line.

In this chapter, you see how to make your way around the curved road of data that leads to nonlinear regression models. The good news is twofold: You can use many of the same techniques you use for regular regression, and in the end, Minitab does the analysis for you.

Anticipating Nonlinear Regression

Nonlinear regression comes into play in situations where you have graphed your data on a *scatterplot* (a two-dimensional graph showing the x variable on the x -axis and the y variable on the y -axis; see the next section "Starting Out with Scatterplots"), and you see a pattern emerging that looks like some type of curve. Examples of data that follow a curve include changes in population size over time, demand for a product as a function of supply, or the length of time that a battery lasts. When a data set follows a curved pattern, the time has come to move away from the linear regression models (covered in

Chapters 4 and 5) and move on to a nonlinear regression model.

Suppose a manager is considering the purchase of new office management software but is hesitating. She wants to know how long it typically takes someone to get up to speed using the software.

What's the statistical question here? She wants a model that shows what the learning curve looks like (on average). (A *learning curve* shows the decrease in time to do a task with more and more practice.) In this scenario, you have two variables: time to complete the task and trial number (for example, the first try is designated by 1, the second try by 2, and so on). Both variables are *quantitative* (numerical) and you want to find a connection between two quantitative variables. At this point, you can start thinking regression.

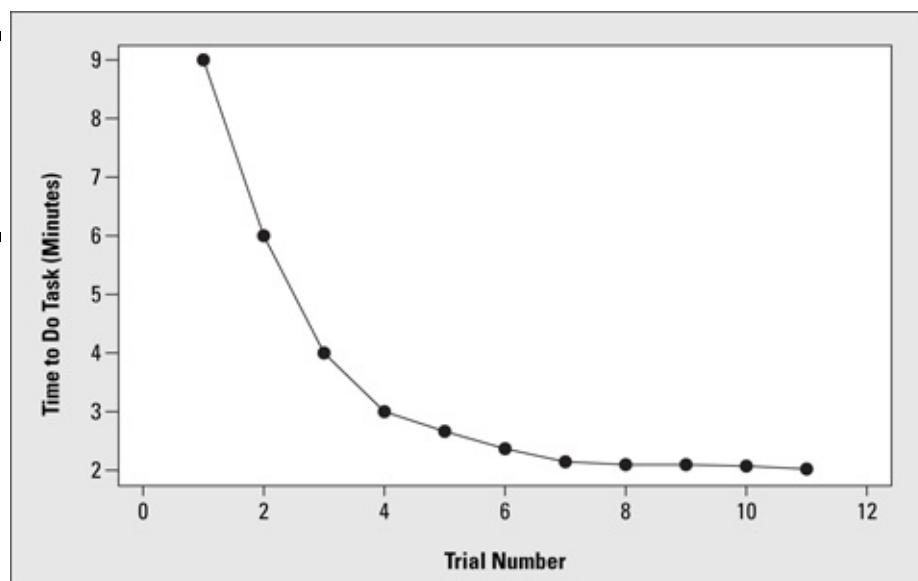
A *regression model* produces a function (be it a line or otherwise) that describes a pattern or relationship. The relationship here is task time versus number of times the task is practiced. But what type of regression model do you use? After all, you can see four types in this book: simple linear regression, multiple regression, nonlinear regression, and logistic regression. You need more clues.

The word “curve” in learning curve is a clue that the relationship being modeled here may not be linear. That word signals that you’re talking about a nonlinear regression model. If you think about what a possible learning curve may look like, you can imagine task time on the *y*-axis and the number of the trial on the *x*-axis.

You may guess that the *y*-values will be high at first, because the first couple of times you try a new task, it takes longer to perform. Then, as the task is repeated, the task time decreases, but at some point more practice doesn’t reduce task time much. So the relationship may be represented by some sort of curve, like the one I simulate in Figure 7-1 (which can be fit by using an exponential function).

This example illustrates the basics of nonlinear regression; the rest of the chapter shows you how the model breaks down.

Figure 7-1:
Learning
curve for time
performing a
new task.



Starting Out with Scatterplots

As with any type of data analysis, before you dive in and select a model that you think fits the data (or that's supposed to fit the data), you have to step back and take a look at the data to see whether any patterns emerge. To do this, look at a scatterplot of the data, and see whether or not you can draw a smooth curve through the data and find that most of the points follow along that curve.

Suppose you're interested in modeling how quickly a rumor spreads. One person knows a secret and tells it to another person, and now two know the secret; each of them tells a person, and now four know the secret; some of those people may pass it on, and so it goes on down the line. Pretty soon, a large number of people know the secret, which is a secret no longer.

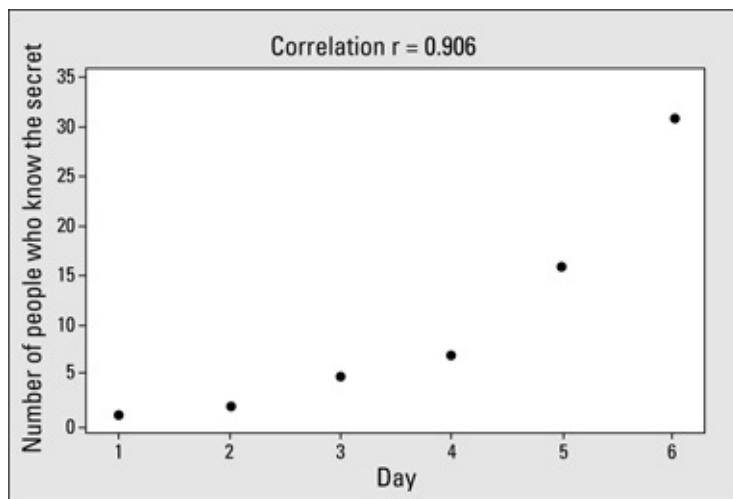
To collect your data, you count the number of people who know a secret by tracking who tells whom over a six-day period. The data are shown in Table 7-1. Note that the spread of the secret catches fire on day 5 — this is how an exponential model works. You can see a scatterplot of the data in Figure 7-2.

Table 7-1 Number of People Knowing a Secret over a 6-Day Period

<i>x</i> (Day)	<i>y</i> (Number of People)
1	1
2	2
3	5
4	7
5	17
6	30

In this situation, the explanatory variable, x , is day, and the response variable, y , is the number of people who know the secret. Looking at Figure 7-2, you can see a pattern between the values of x and y . But this pattern isn't linear. It curves upward. If you tried to fit a line to this data set, how well would it fit?

Figure 7-2: A scatter-plot showing the spread of a secret over a six-day period.



To figure this out, look at the correlation coefficient between x and y , which is found on Figure 7-2 to be 0.906 (see Chapter 4 for more on correlation). You can interpret this correlation as a strong, positive (uphill) linear relationship between x and y . However, in this case, the correlation is misleading because the scatterplot appears to be curved.



If the correlation looks good (close to +1 or -1), don't stop there. As with any regression analysis, it's very important to take into account both the scatterplot and the correlation when making a decision about how well the model being considered would fit the data. The contradiction in this example between the scatterplot and the correlation is a red flag that a straight-line model isn't the best idea.



The correlation coefficient measures only the strength and direction of the linear relationship between x and y (see Chapter 4). However, you may run into situations (like the one shown in Figure 7-2) where a correlation is strong, yet the scatterplot shows a curve would fit better. Don't rely solely on either the scatterplot or the correlation coefficient alone to make your decision about whether to go ahead and fit a straight line to your data.

The bottom line here is that fitting a line to data that appear to have a curved pattern isn't the way to go. Instead, explore models that have curved patterns themselves.

The following sections address two major types of nonlinear (or curved) models that are used to model curved data: polynomials (that are not straight lines — that is, curves like quadratics or cubics), and exponential models (that start out small and quickly increase, or the other way around). Because the pattern of the data in Figure 7-2 starts low and bends upward, the correct model to fit this data is an exponential regression model. (This model is also appropriate for data that start out high and bend down low.)

Handling Curves in the Road with Polynomials

One major family of nonlinear models is the *polynomial* family. You use these models when a polynomial function (beyond a straight line) best describes the curve in the data. (For example, the data may follow the shape of a parabola, which is a second-degree polynomial.) You typically use polynomial models when the data follow a pattern of curves going up and down a certain number of times.

For example, suppose a doctor examines the occurrence of heart problems in patients as it relates to their blood pressure. She finds that patients with very low or very high blood pressure had a higher occurrence of problems, while patients whose blood pressure fell in the middle, constituting the normal range, had fewer problems. This pattern of data has a U-shape, and a para-bola would fit this data well.

In this section, you see what a polynomial regression model is, how you can search for a good-fitting polynomial for your data, and how you can assess polynomial models.

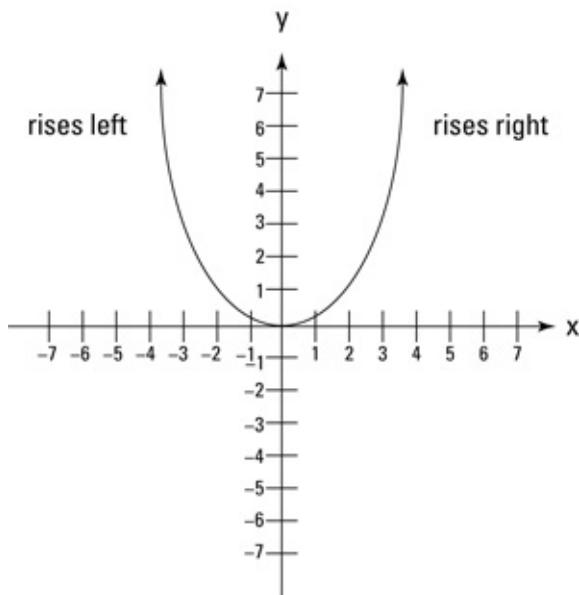
Bringing back polynomials

You may recall from algebra that a *polynomial* is a sum of x terms raised to a variety of powers, and each x is preceded by a constant called the *coefficient* of that term. For example, the model $y = 2x + 3x^2 + 6x^3$ is a polynomial. The general form for a polynomial regression model is $y = \beta_0 + \beta_1x^1 + \beta_2x^2 + \beta_3x^3 + \dots + \beta_kx^k + \varepsilon$. Here, k represents the total number of terms in the model. The ε represents the error that occurs simply due to chance. (Not a bad kind of error, just random fluctuations from a perfect model.)

Here are a few of the more common polynomials you run across when analyzing data and fitting models. Remember, the simplest model that fits is the one you use (don't try to be a hero in statistics — save that for Batman and Robin). The models I discuss in this book are some of your old favorites from algebra: second-, third-, and fourth-degree polynomials.

- ✓ **Second-degree (or quadratic) polynomial:** This model is called a *second-degree (or quadratic) polynomial*, because the largest exponent is 2. An example model is $y = 2x + 3x^2$. A second-degree polynomial forms a parabola shape — either an upside-down or right-side up bowl; it changes direction one time (see Figure 7-3).

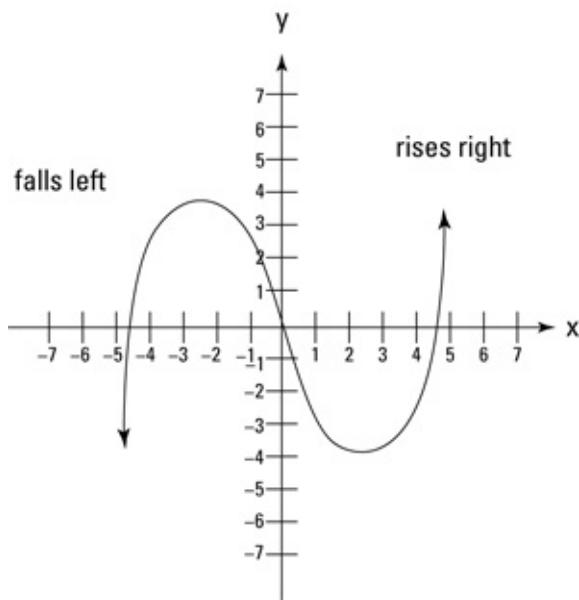
Figure 7-3:
Example of a
second-
degree
polynomial.



- ✓ **Third-degree polynomial:** This model has 3 as the highest power of x . It typically has a sideways S-shape, changing directions two times (see Figure 7-4).
- ✓ **Fourth-degree polynomial:** Fourth-degree polynomials involve x^4 . They typically change directions in curvature three times to look like the letter W or the letter M, depending on whether they're upside down or right-side up (see Figure 7-5).

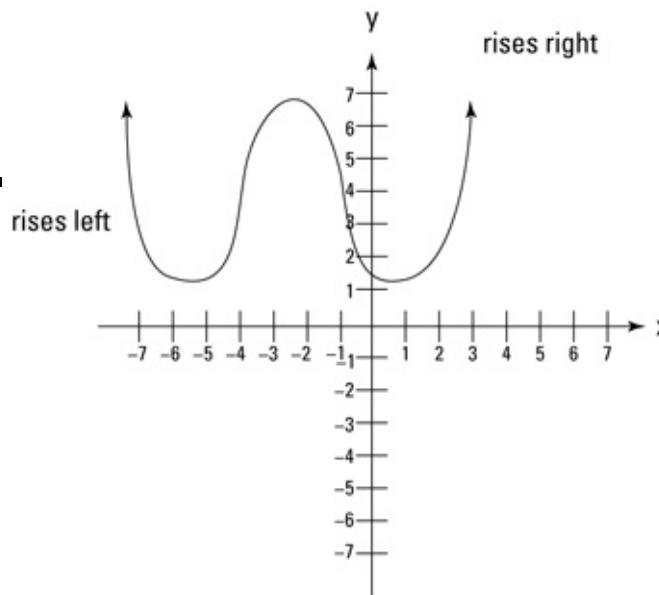
In general, if the largest exponent on the polynomial is n , the number of curve changes in the graph is typically $n - 1$. For more information on graphs of polynomials, refer to your algebra textbook or *Algebra For Dummies* by Mary Jane Sterling (Wiley).

Figure 7-4:
Example of a
third-degree
polynomial.



The nonlinear models in this chapter involve only one explanatory variable, x . You can include more explanatory variables in a nonlinear regression, raising each separate variable to a power, but those models are beyond the scope of this book. I give you information on basic multiple regression models in Chapter 5.

Figure 7-5:
Example of a
fourth-
degree
polynomial.



Searching for the best polynomial model



When fitting a polynomial regression model to your data, you always start with a scatterplot so you can look for patterns; the scatterplot will give you some idea of the type of model that may work. Always start with the simplest model possible and work your way up as needed. Don't plunge in with a high-order polynomial regression model right off the bat. Here are a couple of reasons why:

- ✓ **High-order polynomials are hard to interpret, and their models are complex.** For example, with a straight line, you can interpret the values of the y -intercept and slope easily, but interpreting a tenth-degree polynomial is difficult (and that's putting it mildly).
- ✓ **High-order polynomials tend to cause overfitting.** If you're fitting the model as close as you can to every single point in a data set, your model may not hold for a new data set, meaning that your estimates for y could be way off.



To fit a polynomial to a data set in Minitab, go to Stat>Regression>Fitted Line Plot> and click on the type of regression model you want: linear, quadratic, or cubic. (It doesn't go beyond a third-degree polynomial, but these options should cover 90 percent of the cases where a polynomial is appropriate.) Click on the y variable from the left-hand box, and click Select; this variable will appear in the Response (y) box. Click on the x variable from the left-hand box, and click Select; it will appear in the Predictor (x) box. Click OK.

Following are steps that you can use to see if a polynomial fits your data. (Statistical software can jump in and fit the models for you after you tell it which ones to fit.)

1. Make a scatterplot of your data, and look for any patterns, such as a straight line or a curve.

2. If the data resemble a straight line, try to fit a first-degree polynomial (straight line) to the data first: $y = b_0 + b_1x$.

If the scatterplot doesn't show a linear pattern, or if the correlation isn't close to +1 or -1, move to step three.

3. If the data resemble the shape of a parabola, try to fit a second-degree polynomial: $y = b_0 + b_1x + b_2x^2$.

If the data fit the model well, stop here and refer to the later section "Assessing the fit of a polynomial model." If the model still doesn't fit well, move to step four.

4. If you see curvature that's more complex than a parabola, try to fit a third-degree polynomial: $y = b_0 + b_1x + b_2x^2 + b_3x^3$.

If the data fit the model well, stop here and refer to the later section "Assessing the fit of a polynomial model." If the model still doesn't fit well, move to step five.

5. Continue trying to fit higher-order polynomials until you find one that fits or until the order of the polynomial (largest exponent) simply gets too large to find a reliable pattern.



How large is too large? Typically, if you can't fit the data by the time the degree of the polynomial reaches three, then perhaps a different type of model would work better. Or you may determine that you observe too much scatter and haphazard behavior in the data to try to fit any model.

Minitab can do each of these steps for you up to degree two (that's step three); from there, you need a more sophisticated statistical software program, such as SAS or SPSS. However, most of the models you need to fit go up to the second-degree polynomials.

Using a second-degree polynomial to pass the quiz

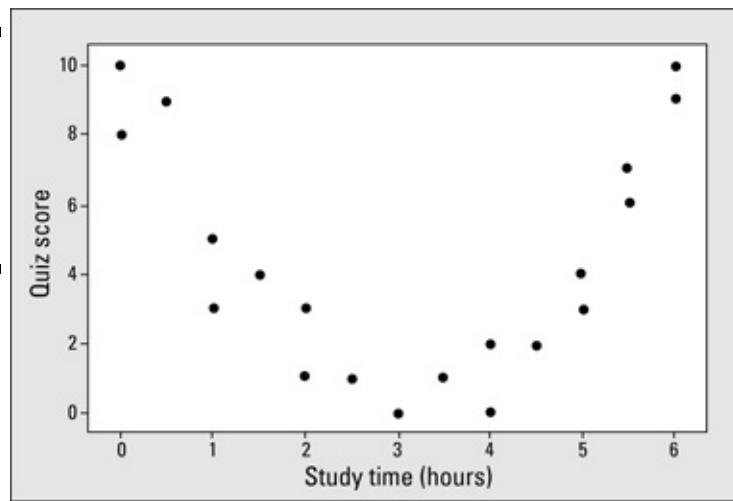
The first step in fitting a polynomial model is to graph the data in a scatterplot and see whether the data fall into a particular pattern. Many different types of polynomials exist to fit data that have a curved type of pattern. One of the most common patterns found in curved data is the quadratic pattern, or second-degree polynomial, which goes up and comes back down, or goes down and comes back up, as the x values move from left to right (see Figure 7-3). The second-degree polynomial is the simplest and most commonly used polynomial beyond the straight line, so it deserves special consideration. (After you master the basic ideas based on second-degree polynomials, you can apply them to polynomials with higher powers.)

Suppose 20 students take a statistics quiz. You record the quiz scores (which have a maximum score of ten) and the number of hours students reported studying for the quiz. You can see the results in Figure 7-6.

Looking at Figure 7-6, it appears that three camps of students are in this class. Camp 1, on the left end of the x -axis, understands the stuff (as reflected in their higher scores) but didn't have to study hardly at all (see that their study time on the x -axis is low). Camp 3 also did very well on the quiz (as indicated by their high quiz scores) but had to study a great deal to get those grades (as seen on the far-right end of the x -axis). The students in the middle, Camp 2, didn't seem to fare well.

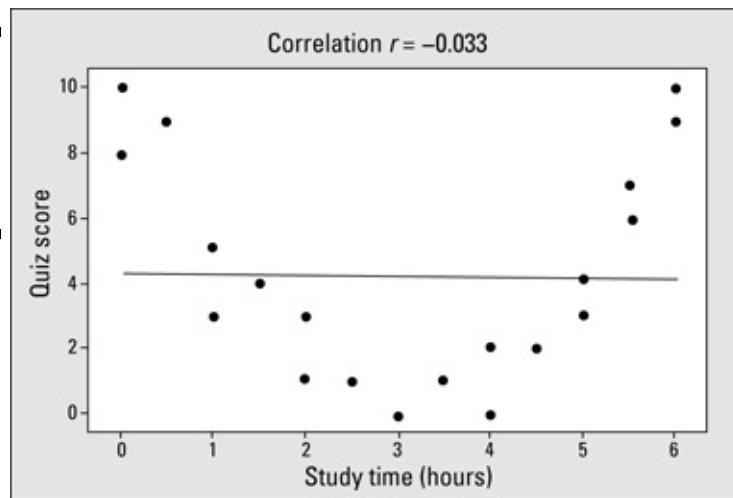
All in all, based on the scatterplot, it does appear that study time may explain quiz scores on some level in a way indicative of a second-degree polynomial. So a quadratic regression model may fit this data.

Figure 7-6:
Scatterplot
showing
study time
and quiz
scores.



Suppose a data analyst (not you!) doesn't know about polynomial regression and just tries to fit a straight line to the quiz-score data. In Figure 7-7, you can see the data and the straight line that he tried to fit to the data. The correlation as shown in the figure is -0.033 , which is basically zero. This correlation means that no linear relationship lies between x and y . It doesn't mean that no relationship is present at all, just that it's not a linear relationship (see Chapter 4 for more on linear relationships). So trying to fit a straight line here was indeed a bad idea.

Figure 7-7:
Trying to fit a
straight line
to quadratic
data.



After you know that a quadratic polynomial seems to be a good fit for the data, the next

challenge is finding the equation for that particular parabola that fits the data from among all the possible parabolas out there.

Remember from algebra that the general equation of a parabola is $y = ax^2 + bx + c$. Now you have to find the values of a , b , and c that create the best-fitting parabola to the data (just like you find the a and the b that create the best-fitting line to data in a linear regression model). That's the object of any regression analysis.

Suppose that you fit a quadratic regression model to the quiz-score data by using Minitab (see the Minitab output in Figure 7-8 and the instructions for using Minitab to fit this model in the previous section). On the top line of the output, you can see that the equation of the best-fitting parabola is quiz score = $9.82 - 6.15 * (\text{study time}) + 1.00 * (\text{study time})^2$. (Note that y is quiz score and x is study time in this example because you're using study time to predict quiz score.)

Figure 7-8: Minitab output for fitting a parabola to the quiz-score data.

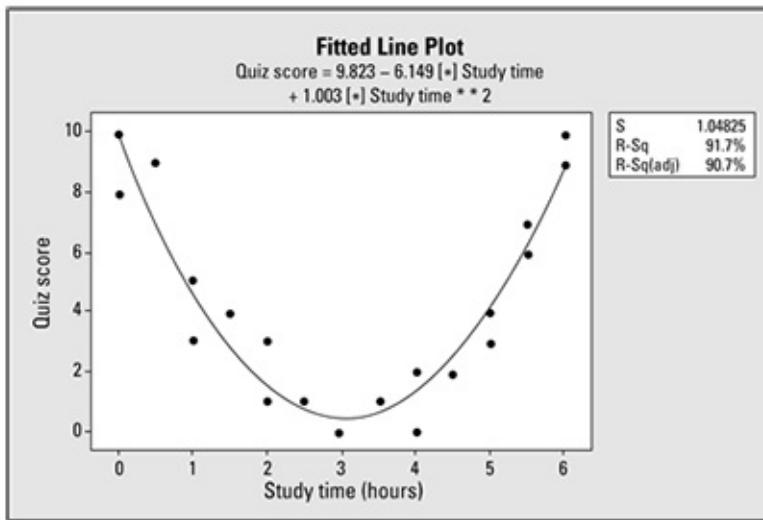
Polynomial Regression Analysis: Quiz Score versus Study Time

The regression equation is
Quiz score = $9.823 - 6.149 \text{ study time} + 1.003 \text{ study time}^{**2}$
 $S = 1.04825 \quad R-Sq = 91.7\% \quad R-Sq(\text{adj}) = 90.7\%$

The scatterplot of the quiz-score data and the parabola that was fit to the data via the regression model is shown in Figure 7-9. From algebra, you may remember that a positive coefficient on the quadratic term (here $a = 1.00$) means the bowl is right-side up, which you can see is the case here.

Looking at Figure 7-9, it appears that the quadratic model fits this data pretty well, because the data fall close to the curve that Minitab found. However, data analysts can't live by scatterplots alone, so the next section helps you figure out how to assess the fit of a polynomial model in more detail.

Figure 7-9:
The parabola appears to fit the quiz-score data nicely.



Assessing the fit of a polynomial model

You make a scatterplot of your data, and you see a curved pattern. So you use polynomial regression to fit a model to the data; the model appears to fit well because the points follow closely to the curve Minitab found, but don't stop there. To make sure your results can be generalized to the population from which your data was taken, you need to do a little more checking beyond just the graph to make sure your model fits well.

To assess the fit of any model beyond the usual suspect, a scatterplot of the data, you look at two additional items, typically in this order: the value of R^2 adjusted and the residual plots.



All three assessments must agree before you can conclude that the model fits. If the three assessments don't agree, you'll likely have to use a different model to fit the data besides a polynomial model, or you'll have to change the units of the data to help a polynomial model fit better. However, the latter fix is outside the scope of Stats II, and you probably won't encounter that situation.

In the following sections, you take a deeper look at the value of R^2 adjusted and the residual plots and figure out how you can use them to assess your model's fit. (You can find more info on the scatterplot in the section "Starting out with Scatterplots" earlier in this chapter.)

Examining R^2 and R^2 adjusted

Finding R^2 , the *coefficient of determination* (see Chapter 5 for full details), is like the day of reckoning for any model. You can find R^2 on your regression output, listed as “R-Sq” right under the portion of the output where the coefficients of the variables appear. Figure 7-8 shows the Minitab output for the quiz-score data example; the value of R^2 in this case is 91.7 percent.

The value of R^2 tells you what percentage of the variation in the y -values the model can explain. To interpret this percentage, note R^2 is the square of r , the correlation coefficient (see Chapter 5). Because values of r beyond ± 0.80 are considered to be good, R^2 values above 0.64 are considered pretty good also, especially for models with only one x variable.

You can consider values of R^2 over 80 percent good, and values under 60 percent aren’t good. Those in between I’d consider so-so; they could be better. (This assessment is just my rule of thumb; opinions may vary a bit from one statistician to another.)

However, you can find such a thing in statistics as too many variables spoiling the pot. Every time you add another x variable to a regression model, the value of R^2 automatically goes up, whether the variable really helps or not (this is just a mathematical fact). Right beside R^2 on the computer output from any regression analysis is the value of R^2 adjusted, which adjusts the value of R^2 down a notch for each variable (and each power of each variable) entered into the model. You can’t just throw a ton of variables into a model whose tiny increments all add up to an acceptable R^2 value without taking a hit for throwing everything in the model but the kitchen sink.



To be on the safe side, you should always use R^2 adjusted to assess the fit of your model, rather than R^2 , especially if you have more than one x variable in your model (or more than one power of an x variable). The values of R^2 and R^2 adjusted are close if you have only a couple of different variables (or powers) in the model, but as the number of variables (or powers) increases, so does the gap between R^2 and R^2 adjusted. In that case, R^2 adjusted is the most fair and consistent coefficient to use to examine model fit.

In the quiz-score example (analysis shown in Figure 7-8), the value of R^2 adjusted is 90.7 percent, which is still a very high value, meaning that the quadratic model fits this data very well. (See Chapter 6 for more on R^2 and R^2 adjusted.)

Checking the residuals

You’ve looked at the scatterplot of your data and the value of R^2 is high. What’s next? Now you examine how well the model fits each individual point in the data to make sure you can’t find any spots where the model is way off or places where you missed another underlying pattern in the data.

A *residual* is the amount of error, or leftover, that occurs when you fit a model to a data set. The residuals are the distances between the predicted values in the model and the observed values of the data themselves. For each observed y -value in the data set, you also have a predicted value from the model, typically called *y-hat*, denoted \hat{y} . The residual is the difference between the values of y and \hat{y} . Each y -value in the data set has a residual; you examine all the residuals together as a group, looking for patterns or unusually high values (indicating a big difference between the observed y and the predicted \hat{y} at that point; see Chapter 4 for the full info on residuals and their plots).

In order for the model to fit well, the residuals need to meet two conditions:

- ✓ **The residuals are independent.** The independence of residuals means that you don't see any pattern as you plot the residuals. The residuals don't affect each other and should be random.
- ✓ **The residuals have a normal distribution centered at zero, and the standardized residuals follow suit.** Having a normal distribution with mean zero means that most of the residuals should be centered around zero, with fewer of them occurring the farther from zero you get. You should observe about as many residuals above the zero line as below it. If the residuals are standardized, this means that as a group their standard deviation is 1; you should expect about 95 percent of them to lie between -2 and +2, following the 68-95-99.7 Rule (see your Stats I text).

You determine whether or not these two conditions are met for the residuals by using a series of four graphs called *residual plots*. Most statisticians prefer to standardize the residuals (meaning they convert them to Z-scores by subtracting their mean and dividing by their standard deviation) before looking at them, because then they can compare the residuals with values on a Z-distribution. If you take this step also, you can ask Minitab to give you a series of four standardized residual plots with which to check the conditions. (See Chapter 4 for full details on standardized residuals and residual plots.)

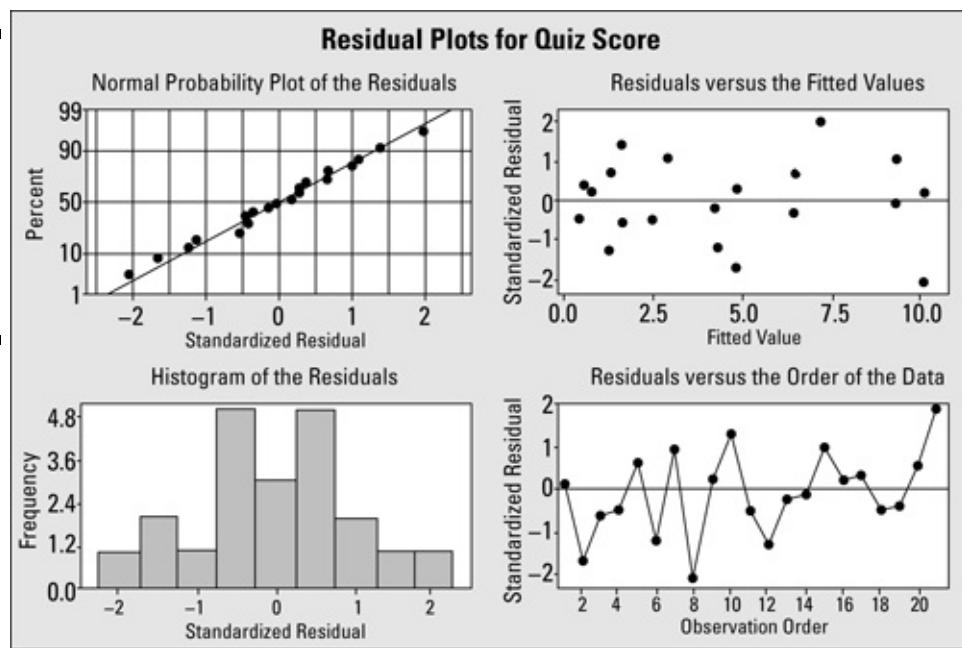
Figure 7-10 shows the standardized residual plots for the quadratic model, using the quiz-score data from the previous sections.

- ✓ The upper-left plot shows that the standardized residuals resemble a normal distribution because your data and the normal distribution match up pretty well, point for point.
- ✓ The upper-right plot shows that most of the standardized residuals fall between -2 and +2 (see Chapter 4 for more on standardized residuals).
- ✓ The lower-left plot shows that the residuals bear some resemblance to a normal distribution.
- ✓ The lower-right plot demonstrates how the residuals have no pattern. They appear to occur at random.

When taken together, all these plots suggest that the conditions on the residual are met

to apply the selected quadratic regression model.

Figure 7-10:
Standardized residual plots for the quiz-score data, using the quadratic model.



Making predictions

After you've found the model that fits well, you can use that model to make predictions for y given x : Simply plug in the desired x -value, and out comes your predicted value for y . (Make sure any values you plug in for x occur within the range of where data were collected; if not, you can't guarantee the model holds.)

Returning to the quiz-score data from previous sections, can you use study time to predict quiz score by using a quadratic regression model? By looking at the scatterplot and the value of R^2 adjusted (review Figures 7-8 and 7-9, respectively), you can see that the quadratic regression model appears to fit the data well. (Isn't it nice when you find something that fits?) The residual plots in Figure 7-10 indicate that the conditions seem to be met to fit this model; you can find no major patterns in the residuals, they appear to center at 1, and most of them stay within the normal boundaries of standardized residuals of -2 and +2.

Considering all this evidence together, study time does appear to have a quadratic relationship with quiz score in this case. You can now use the model to make estimates of quiz score given study time. For example, because the model (shown in Figure 7-8) is $y = 9.82 - 6.15x + 1.00x^2$, if your study time is 5.5 hours, then your estimated quiz score is $9.82 - 6.15 * 5.5 + 1.00 * 5.5^2 = 9.82 - 33.83 + 30.25 = 6.25$. This value makes sense according to what you see on the graph in Figure 7-6 if you look at the place where $x = 5.5$; the y -values are in the vicinity of 6 to 7.



As with any regression model, you can't estimate the value of y for x -values outside the range of where data was collected. If you try to do this, you commit a no-

no called extrapolation. It refers to trying to make predictions beyond where your data allows you to. You can't be sure that the model you fit to your data actually continues ad infinitum for any old value of x . In the quiz-score example (see Figure 7-6), you really can't estimate quiz scores for study times higher than six hours using this model because the data doesn't show anyone studying more than six hours. The model likely levels off after six hours to a score of ten, indicating that studying more than six hours is overkill. (You didn't hear that from me though!)

Going Up? Going Down? Go Exponential!

Exponential models work well in situations where a y variable either increases or decreases exponentially over time. That means the y variable either starts out slow and then increases at a faster and faster rate or starts out high and decreases at a faster and faster rate.

Many processes in the real world behave like an exponential model: for example, the change in population size over time, average household incomes over time, the length of time a product lasts, or the level of patience one has as the number of statistics homework problems goes up.

In this section, you familiarize yourself with the exponential regression model and see how to use it to fit data that either rise or fall at an exponential rate. You also discover how to build and assess exponential regression models in order to make accurate predictions for a response variable y , using an explanatory variable x .

Recollecting exponential models

Exponential models have the form $y = \alpha\beta^x$. These models involve a constant, β , raised to higher and higher powers of x multiplied by a constant, α . The constant β represents the amount of curvature in the model. The constant α is a multiplier in front of the model that shows where the model crosses the y -axis (because when $x = 0$, $y = \alpha * 1$).



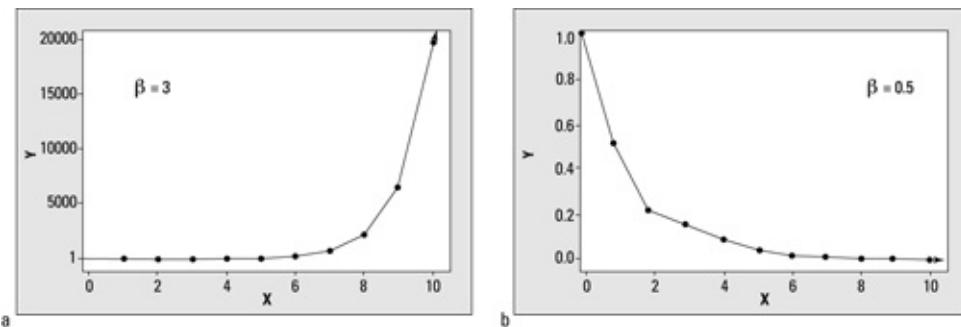
An exponential model generally looks like the upper part of a hyperbola (remember those from advanced algebra?). A hyperbola is a curve that crosses the y -axis at a point and curves downward toward zero or starts at some point and curves upward to infinity (see Figure 7-11 for examples). If β is greater than 1 in an exponential model, the graph curves upward toward infinity. If β is less than 1, the graph curves downward toward zero. All exponential models stay above the x -axis.

For example, the model $y = 1 * 3^x$ is an exponential model. Here, suppose you make $\alpha = 1$, indicating that the model crosses the y -axis at 1 (because plugging $x = 0$ into the equation gives you 1). You set the value of β equal to 3, indicating that you want a bit of curvature to this model. The y -values curve upward quickly from the point $(0, 1)$. For example,

when $x = 1$, you get $1 * 3^1 = 3$; for $x = 2$, you get $1 * 3^2 = 9$; for $x = 3$, you get $1 * 3^3 = 27$; and so on. Figure 7-11a shows a graph of this model. Notice the huge scale needed on the y -axis when x is only 10.

Now suppose you let $\alpha = 1$ and $\beta = 0.5$. These values give you the model $y = 1 * 0.5^x$. This model takes 0.5 (a fraction between 0 and 1) to higher and higher powers starting at $1 * 0.50 = 1$, which makes the y -values smaller and smaller, never reaching zero but always getting closer. (For example, 0.5 to the second power is 0.25, which is less than 0.50, and 0.50 to the tenth power is 0.00098.) Figure 7-11b shows a graph of this model.

Figure 7-11:
The exponential regression model for different values of β .



Searching for the best exponential model

Finding the best-fitting exponential model requires a bit of a twist compared to finding the best-fitting line by using simple linear regression (see Chapter 4). Because fitting a straight-line model is much easier than fitting an exponential model directly from data, you transform the data into something for which a line fits. Then you fit a straight-line model to that transformed data. Finally you undo the transformation, getting you back to an exponential model.

For the transformation, you use *logarithms* because they're the inverse of exponentials. But before you start sweating, don't worry; these math gymnastics aren't something you do by hand — the computer does most of the grunt work for you.

The exponential model looks like this (if you're using base 10): $y = 10^{b_0 + b_1 x}$; note the equation of the line is in the exponent. Follow these steps for fitting an exponential model to your data and using it to make predictions:



The math magic used in these steps is courtesy of the definition of logarithm, which says $\log_b(a) = y \Leftrightarrow b^y = a$. Suppose you have the equation $\log_{10}y = 2 + 3x$.

If you take ten to the power of each side, you get $10^{\log_{10}(y)} = 10^{2+3x}$. By the definition of logarithm, the tens cancel out on the left side and you get $y = 10^{2+3x}$. This model is exponential because x is in the exponent. You can take step two up another notch to include the general form of the straight line model $y = b_0 + b_1 x$. Using the definition of logarithm on this line, you get $\log_{10}(y) = b_0 + b_1 x \Leftrightarrow 10^{b_0 + b_1 x} = y$.

1. Make a scatterplot of the data and see whether the data appears to have a curved pattern that resembles an exponential curve.

If the data follows an exponential curve, proceed to step two; otherwise, consider alternative models (such as multiple regression in Chapter 5).

Chapter 4 tells you how to make a scatterplot in Minitab. For more details on what shape to look for, refer to the previous section.

2. Use Minitab to fit a line to the log(y) data.



In Minitab, you go to the regression model (curve fit). Under Options, select Logten of y. Then select Using scale of logten to give you the proper units for the graph.

Understanding the basic idea of what Minitab does during this step is important; being able to calculate it by hand isn't. Here's what Minitab does:

1. Minitab applies the log (base 10) to the y-values. For example, if y is equal to 100, $\log_{10} 100$ equals 2 (because 10 to the second power equals 100). Note that if the y -values fell close to an exponential model before, the $\log(y)$ values will fall close to a straight-line model. This phenomenon occurs because the logarithm is the inverse of the exponential function, so they basically cancel each other out and leave you with a straight line.

2. Minitab fits a straight line to the $\log(y)$ values by using simple linear regression (see Chapter 4). The equation of the best-fitting straight line for the $\log(y)$ data is $\log(y) = b_0 + b_1x$. Minitab passes this model on to you in its output, and you take it from there.

3. Transform the model back to an exponential model by starting with the straight-line model, $\log(y) = b_0 + b_1x$, that was fit to the $\log_{10}(y)$ data and then applying ten to the power of the left side of the equation and ten to the power of the right side.

By the definition of logarithm, you get y on the left side of the model and ten to the power of $b_0 + b_1x$ on the right side. The resulting exponential model for y is $y = 10^{b_0+b_1x}$.

4. Use the exponential model in step three to make predictions for y (your original variable) by plugging your desired value of x into the model.

Only plug in values for x that are in the range of where the data are located.

5. Assess the fit of the model by looking at the scatterplot of the $\log(y)$ data, checking out the value of R^2 (adjusted) for the straight-line model for $\log(y)$, and checking the residual plots for the $\log(y)$ data.

The techniques and criteria you use to do this are the same as those I discuss in the previous section "Assessing the fit of a polynomial model."

If these steps seem dubious to you, stick with me. The example in the next section lets you see each step firsthand, which helps a great deal. In the end, actually finding predictions by using an exponential model is a lot easier to do than it is to explain.

Spreading secrets at an exponential rate

Often, the best way to figure something out is to see it in action. Using the secret-spreading example from Figure 7-2, you can work through the series of steps from the preceding section to find the best-fitting exponential model and use it to make predictions.

Step one: Check the scatterplot

Your goal in step one is to make a scatterplot of the secret-spreading data and determine whether the data resemble the curved function of an exponential model. Figure 7-2 shows the data for the spread of a number of people knowing the secret, as a function of the number of days. You can see that the number of people who know the secret starts out small, but then as more and more people tell more and more people, the number grows quickly until the secret isn't a secret anymore. This is a good situation for an exponential model, due to the amount of upward curvature in this graph.

Step two: Let Minitab do its thing to $\log(y)$

In step two, you let Minitab find the best-fitting line to the $\log(y)$ data (see the section “Searching for the best exponential model” to find out how to do this in Minitab). The output for the analysis of the secret-spreading data is in Figure 7-12; you can see that the best-fitting line is $\log(y) = -0.19 + 0.28 * x$, where y is the number of people knowing the secret and x is the number of days.

Regression Analysis: Day versus Number

The regression equation is
 $\log_{10}(\text{number}) = -0.1883 + 0.2805 \text{ day}$

S = 0.157335 R-Sq = 93.3% R-Sq(adj) = 91.6%

Figure 7-12:
Minitab fits a line to the $\log(y)$ for the secret-spreading data.

Step three: Go exponential

After you have your Minitab output, you’re ready for step three. You transform the model $\log(y) = -0.19 + 0.28 * x$ into a model for y by taking 10 to the power of the left-hand side and 10 to the power of the right-hand side. Transforming the $\log(y)$ equation for the secret-spreading data, you get $y = 10^{-0.19 + 0.28x}$.

Step four: Make predictions

By using the exponential model from step three, you can move on to step four: Make predictions for appropriate values of x (within the range of where data was collected). Continuing to use the secret-spreading data, suppose you want to estimate the number of people knowing the secret on day five (see Figure 7-2). Just plug $x = 5$ into the exponential model to get $y = 10^{-0.19 + 0.28 * 5} = 10^{1.21} = 16.22$. Looking back at Figure 7-2, you

can see that this estimate falls right in line with the data on the graph.

Step five: Assess the fit of your exponential model

Now that you've found the best-fitting exponential model, you have the worst behind you. You've arrived at step five and are ready to further assess the model fit (beyond the scatterplot of the original data) to make sure no major problems arise.

In general, to assess the fit of an exponential model, you're really looking at the straight-line fit of $\log(y)$. Just use these three items (in any order) in the same way as described in the earlier section "Assessing the fit of a polynomial model":

- ✓ **Check the scatterplot of the $\log(y)$ data to see how well it resembles a straight line.** You assess the fit of the $\log(y)$ for the secret-spreading data first through the scatterplot shown in Figure 7-13. The scatterplot shows that the model appears to fit the data well, because the points are scattered in a tight pattern around a straight line.



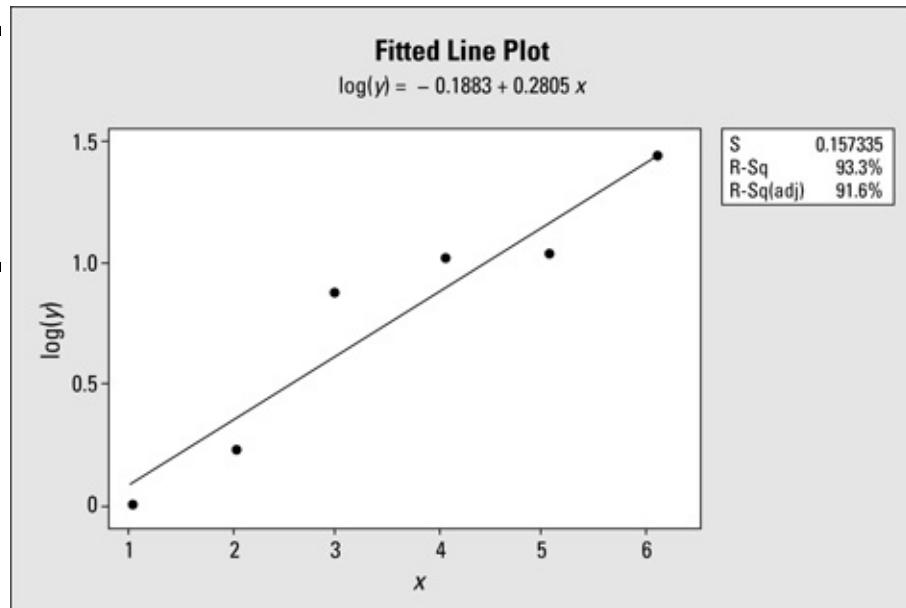
During this process the data were transformed also. You started with x and y data, and now you have x and $\log(y)$ for your data. You see x , y , and $\log(y)$ for the secret-spreading data in Table 7-2.

- ✓ **Examine the value of R^2 adjusted for the model of the best-fitting line for $\log(y)$, done by Minitab.** The value of R^2 adjusted for this model is found in Figure 7-13 to be 91.6 percent. This value also indicates a good fit because it's very close to 100 percent. Therefore, 91.6 percent of the variation in the number of people knowing the secret is explained by how many days it has been since the secret-spreading started. (Makes sense.)
- ✓ **Look at the residual plots from the fit of a line to the $\log(y)$ data.** The residual plots from this analysis (see Figure 7-14) show no major departures from the conditions that the errors are independent and have a normal distribution. Note that the histogram in the lower-left corner doesn't look all that bell-shaped, but you don't have a lot of data in this example, and the rest of the residual plots seem okay. So, you have little cause to really worry.

Table 7-2 Log(y) Values for the Secret-Spreading Data

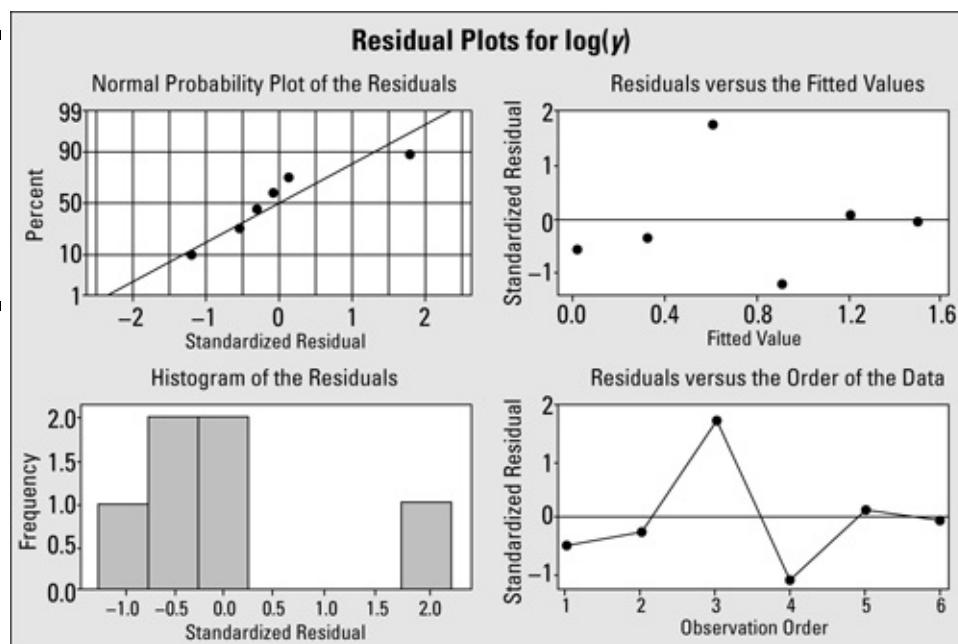
x (Day)	y (Number of People)	$\log(y)$
1	1	0.00
2	2	0.30
3	5	0.70
4	7	0.85
5	17	1.23

Figure 7-13: A scatter-plot showing the fit of a straight line to $\log(y)$ data.



All in all, it appears that the secret's out on the secret-spreading data, now that you have an exponential model that explains how it happens.

Figure 7-14: Residual plots showing the fit of a straight line to $\log(y)$ data.



Chapter 8

Yes, No, Maybe So: Making Predictions by Using Logistic Regression

In This Chapter

- ▶ Knowing when logistic regression is appropriate
 - ▶ Building logistic regression models for yes or no data
 - ▶ Checking model conditions and making the right conclusions
-

Everyone (even yours truly) tries to make predictions about whether or not a certain event is going to happen. For example, what's the chance it's going to rain this weekend? What are your team's chances of winning the next game? What's the chance that I'll have complications during this surgery? These predictions are often based on *probability*, the long-term percentage of time an event is expected to happen.

In the end, you want to estimate p , the probability of an event occurring. In this chapter, you see how to build and test models for p based on a set of explanatory (x) variables. This technique is called *logistic regression*, and in this chapter, I explain how to put it to use.

Understanding a Logistic Regression Model

In a logistic regression, you're estimating the probability that an event occurs for a randomly selected individual versus the probability that the event doesn't occur. In essence, you're looking at yes or no data: yes it occurred (probability = p); or no, it didn't occur (probability = $1 - p$). Yes or no data that come from a random sample have a binomial distribution with probability of success (the event occurring) equal to p .

In the binomial problems you saw in Stats I, you had a sample of size n trials, you had yes or no data, and you had a probability of success on each trial, denoted by p . In your Stats I course, for any binomial problem the value of p was somehow given to be a certain value, like a fair coin has probability $p = 0.50$ for coming up heads. But in Stats II, you operate under the much more realistic scenario that it's not. In fact, because p isn't known, your job is to estimate what it is and use a model to do that.



To estimate p , the chance of an event occurring, you need data that come in the form of yes or no, indicating whether or not the event occurred for each individual in the data set.

Because yes or no data don't have a normal distribution, which is a condition needed for other types of regression, you need a new type of regression model to do this job; that model is *logistic regression*.

How is logistic regression different from other regressions?

You use logistic regression when you use a quantitative variable to predict or guess the outcome of some categorical variable with only two outcomes (for example, using barometric pressure to predict whether or not it will rain).

A logistic regression model ultimately gives you an estimate for p , the probability that a particular outcome will occur in a yes or no situation (for example, the chance that it will rain versus not). The estimate is based on information from one or more explanatory variables; you can call them $x_1, x_2, x_3, \dots, x_k$. (For example, x_1 = humidity, x_2 = barometric pressure, x_3 = cloud cover, . . . and x_k = wind speed.)

Because you're trying to use one variable (x) to make a prediction for another variable (y), you may think about using regression — and you would be right. However, you have many types of regression to choose from, and you need to determine what kind is most appropriate here. You need the type of regression that uses a quantitative variable (x) to predict the outcome of some categorical variable (y) that has only two outcomes (yes or no).

So being the good Stats II student that you are, you go to your trusty list of statistical techniques, and you look under regression — and immediately see more than one type.

- ✓ You see simple linear regression. No, you use that when you have one quantitative variable predicting another (see Chapter 4).
- ✓ Multiple regression? No, that method just expands simple linear regression to add more x variables (see Chapter 5).
- ✓ Nonlinear regression? Well no, that still works with two quantitative variables; it's just that the data form a curve, not a line.

But then you come across logistic regression, and . . . eureka! You see that logistic regression handles situations where the x variable is numerical and the y variable is categorical with two possible categories. Just what you're looking for!



Logistic regression, in essence, estimates the probability of y being in one category or the other, based on the value of some quantitative variable, x . For example, suppose you want to predict someone's height based on gender. Because gender is a categorical variable, you use logistic regression to make these

predictions. Suppose a 1 indicates a male. People who receive a probability of more than 0.5 of being male (based on their heights) are predicted to be male, and people who receive a probability of less than 0.5 of being male (based on their heights) are predicted to be female.



In this chapter, I present only the case where you use one explanatory variable to predict the outcome. You can extend the ideas in exactly the same way as you can extend the simple linear regression model to a multiple regression model.

Using an S-curve to estimate probabilities

In a simple linear regression model, the general form of a straight line is $y = \beta_0 + \beta_1x$ and y is a quantitative variable. In the logistic regression model, the y variable is categorical, not quantitative. What you're estimating, however, is not which category the individual lies in, but rather what the probability is that the individual lies in a certain category. So, the model for logistic regression is based on estimating this probability, called p .

If you were to estimate p using a simple linear regression model, you may think you should try to fit a straight line, $p = \beta_0 + \beta_1x$. However, it doesn't make sense to use a straight line to estimate the probability of an event occurring based on another variable, due to the following reasons:

- ✓ **The estimated values of p can never be outside of [0, 1], which goes against the idea of a straight line (a straight line continues on in both directions).**
- ✓ **It doesn't make sense to force the values of p to increase in a linear way based on x .** For example, an event may occur very frequently with a range of large values of x and very frequently with a range of small values of x , with very little chance of the event happening in an area in between. This type of model would have a U shape rather than a straight-line shape.

To come up with a more appropriate model for p , statisticians created a new function of p whose graph is called an S-curve. The *S-curve* is a function that involves p , but it also involves e (the natural logarithm) as well as a ratio of two functions.

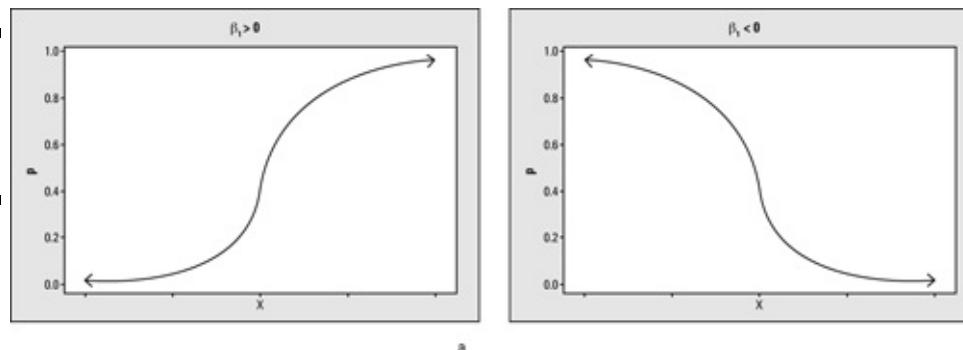
The values of the S-curve always fit between 0 and 1, which allows the probability, p , to change from low to high or high to low, according to a curve that's shaped like an S. The general form of the logistic regression model based on an S-curve is
$$p = \frac{e^{\beta_0 + \beta_1x}}{1 + e^{\beta_0 + \beta_1x}}$$
.

Interpreting the coefficients of the logistic regression model



The sign on the parameter β_1 tells you the direction of the S-curve. If β_1 is positive, the S-curve goes from low to high (see Figure 8-1a); if β_1 is negative, the S-curve goes from high to low (Figure 8-1b).

Figure 8-1:
Two basic types of S-curves.



The magnitude of β_1 (indicated by its absolute value) tells you how much curvature is in the model. High values indicate a steep curvature, and low values indicate gradual curvature. The parameter β_0 just shifts the S-curve to the proper location to fit your data. It shows you the cutoff point where x -values change from high to low probability and vice versa.

The logistic regression model in action

Often, the best way to figure something out is to see it in action. In this section, I give you an example of a situation where you can use a logistic regression model to estimate a probability. (I expand on this example later in this chapter; for now, I'm just setting up a scenario for logistic regression.)

Suppose movie marketers want to estimate the chance that someone will enjoy a certain family movie, and you believe age may have something to do with it. Translating this research question into x 's and y 's, the response variable (y) is whether or not a person will enjoy the movie, and the explanatory variable (x) is the person's age. You want to estimate p , the chance of someone enjoying the movie.

You collect data on a random sample of 40 people, shown in Table 8-1. Based on your data, it appears that younger people enjoyed the movie more than older people and that at a certain age, the trend switches from liking the movie to disliking it. Armed with this data, you can build a logistic regression model to estimate p .

Table 8-1 Movie Enjoyment (Yes or No Data) Based on Age

Age	Enjoyed the Movie	Total Number Sampled
10	3	3
15	4	4
16	3	3

18	2	3
20	2	3
25	2	4
30	2	4
35	1	5
40	1	6
45	0	3
50	0	2

Carrying Out a Logistic Regression Analysis

The basic idea of any model-fitting process is to look at all possible models you can have under the general format and find the one that fits your data best.

The general form of the best-fitting logistic regression model is $\hat{p} = \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}}$, where \hat{p} is the estimate of p , b_0 is the estimate of β_0 , and b_1 is the estimate of β_1 (from the previous section “Using an S-curve to estimate probabilities”). The only values you have a choice about to form your particular model are the values of b_0 and b_1 . These values are the ones you’re trying to estimate through the logistic regression analysis.

To find the best-fitting logistic regression model for your data, complete the following steps:

- 1. Run a logistic regression analysis on the data you collected (see the next section).**
- 2. Find the coefficients of constant and x , where x is the name of your explanatory variable.**

These coefficients are b_0 and b_1 , the estimates of β_0 and β_1 in the logistic regression model.

- 3. Plug the coefficients from step one into the logistic regression model:**

$$\hat{p} = \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}}.$$

This equation is your best-fitting logistic regression model for the data. Its graph is an S-curve (for more on the S-curve, see the section “Using an S-curve to estimate probabilities” earlier in this chapter).

In the sections that follow, you see how to ask Minitab to do the above steps for you. You also see how to interpret the resulting computer output, find the equation of the best-fitting logistic regression model, and use that model to make predictions (being ever mindful that all conditions are met).

Running the analysis in Minitab



Here's how to perform a logistic regression using Minitab (other statistical software packages are similar):

1. Input your data in the spreadsheet as a table that lists each value of the x variable in column one, the number of yeses for that value of x in column two, and the total number of trials at that x -value in column three.

These last two columns represent the outcome of the response variable y . (For an example of how to enter your data, see Table 8-1 based on the movie and age data.)

2. Go to Stat>Regression>Binary Logistic Regression.

3. Beside the Success option, select your variable name from column two, and beside Trial, select your variable name for column three.

4. Under Model, select your variable name from column one, because that's the column containing the explanatory (x) variable in your model.

5. Click OK, and you get your logistic regression output.

When you fit a logistic regression model to your data, the computer output is composed of two major portions:

- ✓ **The model-building portion:** In this part of the output, you can find the coefficients b_0 and b_1 . (I describe coefficients in the next section.)
- ✓ **The model-fitting portion:** You can see the results of a Chi-square goodness-of-fit test (see Chapter 15) as well as the percentage of concordant and discordant pairs in this section of the output. (A *concordant pair* means the predicted outcome from the model matches the observed outcome from the data. A *discordant pair* is one that doesn't match.)

In the case of the movie and age data, the model-building part of the Minitab output is shown in Figure 8-2. The model-fitting part of the Minitab output from the logistic regression analysis is in Figure 8-4.

In the following sections, you see how to use this output to build the best-fitting logistic regression model for your data and to check the model's fit.

Figure 8-2:

The model-building part of the movie and age data's logistic regression output.

Predictor	Coef	SE Coef	Z	P	Odds Ratio	95% CI Lower	95% CI Upper
Constant	4.86539	1.43434	3.39	0.001			
Age	-0.175745	0.0499620	-3.52	0.000	0.84	0.76	0.93

Finding the coefficients and making the model

After you have Minitab run a logistic regression analysis on your data, you can find the coefficients b_0 and b_1 and put them together to form the best-fitting logistic regression model for your data.

Figure 8-2 shows part of the Minitab output for the movie enjoyment and age data. (I discuss the remaining output in the section “Checking the fit of the model.”) The first column of numbers is labeled Coef, which stands for the coefficients in the model. The first coefficient, b_0 , is labeled Constant. The second coefficient is in the row labeled by your explanatory variable, x . (In the movie and age data, the explanatory variable is age. This age coefficient represents the value of b_1 in the model.)

According to the Minitab output in Figure 8-2, the value of b_0 is 4.87 and the value of b_1 is -0.18. After you’ve determined the coefficients b_0 and b_1 from the Minitab output to find the best-fitting S-curve for your data, you put these two coefficients into the general logistic regression model: $\hat{p} = \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}}$.

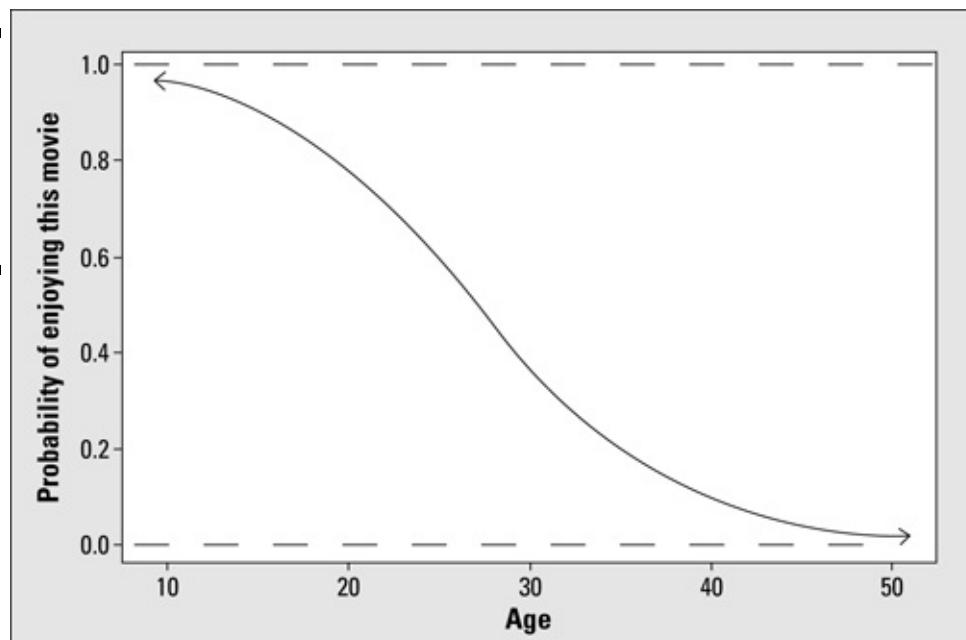
For the movie and age data, you get $\hat{p} = \frac{e^{4.87 - 0.18x}}{1 + e^{4.87 - 0.18x}}$, which is the best-fitting logistic regression model for this data set.

The graph of the best-fitting logistic regression model for the movie and age data is shown in Figure 8-3. Note that the graph is a downward-sloping S-curve because higher probabilities of liking the movie are affiliated with lower ages and lower probabilities are affiliated with higher ages.

The movie marketers now have the answer to their question. This movie has a higher chance of being well liked by kids (and the younger, the better) and a lower chance of being well liked by adults (and the older they are, the lower the chance of liking the movie).

The point where the probability changes from high to low (that is, at the $\hat{p} = 0.50$ mark) is between ages 25 and 30. That means that the tide of probability of liking the movie appears to turn from higher to lower in that age range. Using calculus terms, this point is called the *saddle point* of the S-curve, which is the point where the graph changes from concave up to concave down, or vice versa.

Figure 8-3:
The best-fitting S-curve for the movie and age data.



Estimating p

You've determined the best-fitting logistic regression model for your data, obtained the values of b_0 and b_1 from the logistic regression analysis, and know the precise S-curve that fits your data best (check out the previous sections). You're now ready to estimate p and make predictions about the probability that the event of interest will happen, given the value of the explanatory variable x .

To estimate p for a particular value of x , plug that value of x into your equation (the best-fitting logistic regression model) and simplify it by using your algebra skills. The number you get is the estimated chance of the event occurring for that value of x , and it should be a number between 0 and 1, being a probability and all.

Continuing with the movie and age example from the preceding sections, suppose you want to predict whether a 15-year-old would enjoy the movie. To estimate p , plug 15 in for

$$x \text{ in the logistic regression model } \hat{p} = \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}} \text{ to get } \hat{p} = \frac{e^{4.87 - 0.18 \cdot 15}}{1 + e^{4.87 - 0.18 \cdot 15}} = \frac{e^{2.17}}{1 + e^{2.17}} = \frac{8.76}{9.76} = 0.90.$$

That answer means you estimate there's a 90 percent chance that a 15-year-old will like the movie. You can see in Figure 8-3 that when x is 15, p is approximately 0.90. On the other hand, if the person is 50 years old, the chance he'll like this movie is $\hat{p} = \frac{e^{4.87 - 0.18 \cdot 50}}{1 + e^{4.87 - 0.18 \cdot 50}}$, or 0.02 (shown in Figure 8-3 for $x = 50$), which is only a 2 percent chance.

Checking the fit of the model



The results you get from a logistic regression analysis, as with any other data analysis, are all subject to the model fitting appropriately.

To determine whether or not your logistic regression model fits, follow these steps (which I cover in more detail later in this section):

1. Locate the p -value of the goodness-of-fit test (found in the Goodness-of-Fit portion of the computer output; see Figure 8-4 for an example). If the p -value is larger than 0.05, conclude that your model fits, and if the p -value is less than 0.05, conclude that your model doesn't fit.
2. Find the p -value for the b_1 coefficient (it's listed under P in the row for your column one (explanatory) variable in the model-building portion of the output; see Figure 8-2 for an example). If the p -value is less than 0.05, the x variable is statistically significant in the model, so it should be included. If the p -value is greater than or equal to 0.05, the x variable isn't statistically significant and shouldn't be included in the model.
3. Look later in the output at the percentage of concordant pairs. This percentage

reflects the proportion of time that the data and the model actually agree with each other. The higher the percentage, the better the model fits.



The conclusion in step one based on the p -value may seem backward to you, but here's what's happening: Chi-square goodness-of-fit tests measure the overall difference between what you expect to see via your model and what you actually observe in your data. (Chapter 15 gives you the lowdown on Chi-square tests.) The null hypothesis (H_0) for this test says you have a difference of zero between what you observed and what you expected from the model; that is, your model fits. The alternative hypothesis, denoted H_a , says that the model doesn't fit. If you get a small p -value (under 0.05), reject H_0 and conclude the model doesn't fit. If you get a larger p -value (above 0.05), you can stay with your model.



Failure to reject H_0 here (having a large p -value) only means that you can't say your model doesn't fit the population from which the sample came. It doesn't necessarily mean the model fits perfectly. Your data could be unrepresentative of the population just by chance.

Fitting the movie model

You're ready to check out the fit of the movie data to make sure you still have a job when the box office totals come in.

Step one: p -value for Chi-squared

Using Figure 8-4 to complete the first step of checking the model's fit, you can see many different goodness-of-fit tests. The particulars of each of these tests are beyond the scope of this book; however, in this case (as with most cases), each test has only slightly different numerical results and the same conclusions.

All the p -values in column four of Figure 8-4 are over 0.80, which is much higher than the 0.05 you need to reject the model. After looking at the p -values, the model using age to predict movie likeability appears to fit this data.

Goodness-of-Fit Test			
Method	Chi-Square	DF	P
Pearson	2.83474	9	0.970
Deviance	3.63590	9	0.934
Hosmer-Lemeshow	2.75232	6	0.839

Measures of Association:			
(Between the Response Variable and Predicted Probabilities)			
Pairs	Number	Percent	Summary Measures
Concordant	349	87.3	Somers' D 0.80
Discordant	30	7.5	Goodman-Kruskal Gamma 0.84
Ties	21	5.3	Kendall's Tau-a 0.41
Total	400	100.0	

Figure 8-4:
The model-fitting part of the movie and age data's logistic regression output.

Step two: p-value for the x variable

For step two, you look at the significance of the x variable age. Back in Figure 8-2, you can see the constant for age, -0.18 , and farther along in its row, you can see that the Z-value is -3.52 ; this Z-value is the test statistic for testing $H_0: \beta_1 = 0$ versus $H_a: \beta_1 \neq 0$. The p -value is listed as 0.000 , which means it's smaller than 0.001 (a highly significant number). So you know that the coefficient in front of x , also known as β_1 , is statistically significant (not equal to zero), and you should include x (age) in the model.

Step three: Concordant pairs

To complete step three of the fit-checking process, look at the percentage of concordant pairs reported in Figure 8-4. This value shows the percentage of times the data actually agreed with the model (87.3). To determine concordance, the computer makes predictions as to whether the event should have occurred for each individual based on the model and compares those results to what actually happened.



The logistic regression model is for p , the probability of the event occurring, so if p is estimated to be > 0.50 for some value of x , the computer predicts that the event will occur (versus not occurring). If the estimated value of p is < 0.50 for a particular x -value, the computer predicts that it won't occur.

For the movie and age data, the percentage of concordant pairs (that is, the percentage of times the model made the right decision in predicting what would happen) is 87.3 percent, which is quite high.

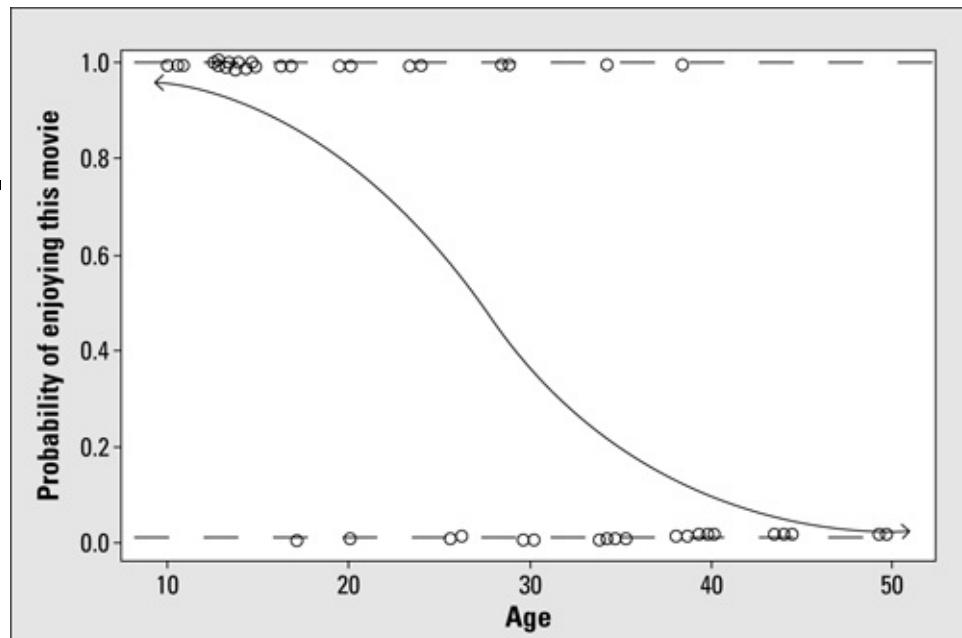


The percentage of concordant pairs was obtained by taking the number of concordant pairs and dividing by the total number of pairs. I'd start getting excited if the percentage of concordant pairs got over 75 percent; the higher, the better.

Figure 8-5 shows the logistic regression model for the movie and age data, with the actual values of the observed data added as circles. The S-curve shows the probability of liking the movie for each age level, and the computer will predict “1” = they will like the movie, if $\hat{p} > 0.50$. Circles indicate whether the people of those age levels actually liked the movie ($y = 1$) or not ($y = 0$).

Much of the time, the model made the right decision; probabilities above 0.50 are associated with more circles at the value of 1, and probabilities below 0.50 are associated with more circles at the value of zero. It's the outcomes that have p near 0.50 that are hard to predict because the results can go either way.

Actual observed values (0 and 1) compared to the model.



Which method to use to compare? Sorting out similar situations

Data come in a variety of forms, and each form has its own analysis to use to make comparisons. It can get difficult to decide which type of analysis to use when.

It may help to sort out some situations that sound similar but have subtle differences that lead to very different analyses. You can use the following list to compare these subtle, but important, differences:

- ✓ If you want to compare three or more groups of numerical variables, use ANOVA (see Chapter 10). For only two groups, use a *t*-test (see Chapters 3 and 9).
 - ✓ If you want to estimate one numerical variable based on another, use simple linear regression (see Chapter 4).
 - ✓ If you want to estimate one numerical variable using many other numerical variables, use multiple regression (see Chapter 5).
 - ✓ If you want to estimate a categorical variable with two categories by using a numerical variable, use logistic regression, which is the focus of this chapter, of course.
 - ✓ If you want to compare two categorical variables to each other, head straight for a Chi-square test (see Chapter 14).

All this evidence helps confirm that your model fits your data well. You can go ahead and make predictions based on this model for the next individual that comes up, whose outcome you don't know (see the section “Estimating p ” earlier in this chapter).

Part III

Analyzing Variance with ANOVA

The 5th Wave

By Rich Tennant

© RICH TENNANT



"This is my old Statistics II professor
his wife Doris, and their two children,
Wilcox and Kruskal."

In this part . . .

You get all the nuts and bolts you need to understand one-way and two-way analyses of variance (also known as ANOVA), which compare the means of several populations at one time based on one or two different characteristics. You see how to read and understand ANOVA tables and computer output, and you get to go behind the scenes to understand the big ideas behind the formulas used in ANOVA. Finally, you see lots of different multiple comparison procedures to zoom in on which means are different. (Don't sweat it. I always present formulas only on a need-to-know basis.)

Chapter 9

Testing Lots of Means? Come On Over to ANOVA!

In This Chapter

- ▶ Extending the t -test for comparing two means by using ANOVA
 - ▶ Utilizing the ANOVA process
 - ▶ Carrying out an F -test
 - ▶ Navigating the ANOVA table
-

One of the most commonly used statistical techniques at the Stats II level is *analysis of variance* (affectionately known as ANOVA). Because the name has the word *variance* in it, you may think that this technique has something to do with variance — and you would be right. Analysis of variance is all about examining the amount of variability in a y (response) variable and trying to understand where that variability is coming from.

One way you can use ANOVA is to compare several populations regarding some quantitative variable, y . The populations you want to compare constitute different groups (denoted by an x variable), such as political affiliations, age groups, or different brands of a product. ANOVA is also particularly suitable for situations involving an experiment in which you apply certain treatments (x) to subjects and measure a response (y).

In this chapter, you start with the t -test for two population means, the precursor to ANOVA. Then you move on to the basic concepts of ANOVA to compare more than two means: sums of squares, the F -test, and the ANOVA table. You apply these basics to the *one-factor* or *one-way* ANOVA, where you compare the responses based only on one treatment variable. (In Chapter 11, you can see the basics applied to a two-way ANOVA, which has two treatment variables.)

Comparing Two Means with a t -Test

The *two-sample t -test* is designed to test to see whether two population means are different. The conditions for the two-sample t -test are the following:

- ✓ The two populations are independent. In other words, their outcomes don't affect each other.
- ✓ The response variable (y) is a quantitative variable, meaning its values have numerical meaning and represent quantities of some kind.
- ✓ The y -values for each population have a normal distribution. However, their means may be different; that's what the t -test determines.

✓ The variances of the two normal distributions are equal.



For large sample sizes when you know the variances, you use a Z-test for the two population means. However, a *t*-test allows you to test two population means when the variances are unknown or the sample sizes are small. This occurs quite often in situations where an experiment is performed and the number of subjects is limited. (See your Stats I text or Statistics For Dummies (Wiley) for info on the Z-test.)

Although you've seen *t*-tests before in your Stats I class, it may be good to review the main ideas. The *t*-test tests the hypotheses $H_0: \mu_1 = \mu_2$ versus $H_a: \mu_1 \neq \mu_2$, where the situation dictates which of these hypotheses you use. (**Note:** With ANOVA, you extend this idea to k different means from k different populations, and the only version of H_a of interest is \neq .)



To conduct the two-sample *t*-test, you collect two data sets from the two populations, using two independent samples. To form the test statistic (the *t*-statistic), you subtract the two sample means and divide by the *standard error* (a combination of the two standard deviations from the two samples and their sample sizes). You compare the *t*-statistic to the *t*-distribution with $n_1 + n_2 - 2$ degrees of freedom and find the *p*-value. If the *p*-value is less than the predetermined α level, say 0.05, you have enough evidence to say the population means are different. (For information on hypothesis tests, see Chapter 3.)

For example, suppose you're at a watermelon seed-spitting contest where contestants each put watermelon seeds in their mouths and spit them as far as they can. Results are measured in inches and are treated with the reverence of the shot-put results at the Olympics. You want to compare the watermelon seed-spitting distances of female and male adults. Your data set includes ten people from each group.

You can see the results of the *t*-test in Figure 9-1. The mean spitting distance for females was 47.8 inches; the mean for males was 56.5 inches; and the difference (females – males) is –8.71 inches, meaning the females in the sample spit seeds at shorter distances, on average, than the males. The *t*-statistic for the difference in the two means (females – males) is $t = -2.23$, which has a *p*-value of 0.039 (see the last line of the output in Figure 9-1). At a level of $\alpha = 0.05$, this difference is significant (because $0.039 < 0.05$). You conclude that males and females differ with respect to their mean watermelon seed-spitting distance. And you can say males are likely spitting farther because their sample mean was higher.

mean
watermelon
seed-spitting
distances for
females
versus males.

Two-sample T for females vs males				
	N	Mean	StDev	SE Mean
females	10	47.80	9.02	2.9
males	10	56.50	8.45	2.7
Difference = mu (females) - mu (males)				
Estimate for difference: -8.70000				
95% CI for difference: (-16.90914, -0.49086)				
T-Test of difference = 0 (vs not =): T-Value = -2.23 P-Value = 0.039 DF = 18				

Evaluating More Means with ANOVA

When you can compare two independent populations inside and out, at some point two populations will not be enough. Suppose you want to compare more than two populations regarding some response variable (y). This idea kicks the t -test up a notch into the territory of ANOVA. The ANOVA procedure is built around a hypothesis test called the F -test, which compares how much the groups differ from each other compared to how much variability is within each group. In this section, I set up an example of when to use ANOVA and show you the steps involved in the ANOVA process. You can then apply the ANOVA steps to the following example throughout the rest of the chapter.

Spitting seeds: A situation just waiting for ANOVA

Before you can jump into using ANOVA, you must figure out what question you want answered and collect the necessary data.

Suppose you want to compare the watermelon seed-spitting distances for four different age groups: 6–8 years old, 9–11, 12–14, and 15–17. The hypotheses for this example are $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$ versus $H_a: \text{At least two of these means are different}$, where the population means μ represent those from the age groups, respectively.

Over the years of this contest, you've collected data on 200 children from each age group, so you have some prior ideas about what the distances typically look like. This year, you have 20 entrants, 5 in each age group. You can see the data from this year, in inches, in Table 9-1.

Table 9-1 Watermelon Seed-Spitting Distances for Four Age Groups of Children (Measured in Inches)

6–8 Years	9–11 Years	12–14 Years	15–17 Years
38	38	44	44
39	39	43	47
42	40	40	45
40	44	44	45
41	43	45	46

Do you see a difference in distances for these age groups based on this data? If you were to just combine all the data, you would see quite a bit of difference (the range of the combined data goes from 38 inches to 47 inches). And you may suspect that older kids

can spit farther.

Perhaps accounting for which age group each contestant is in does explain at least some of what's going on. But don't stop there. The next section walks you through the official steps you need to perform to answer your question.

Walking through the steps of ANOVA

You've decided on the quantitative response variable (y) you want to compare for your k various population (or treatment) means, and you've collected a random sample of data from each population (refer to the preceding section). Now you're ready to conduct ANOVA on your data to see whether the population means are different for your response variable, y .

The characteristic that distinguishes these populations is called the *treatment variable*, x . Statisticians use the word *treatment* in this context because one of the biggest uses of ANOVA is for designed experiments where subjects are randomly assigned to treatments, and the responses are compared for the various treatment groups. So statisticians often use the word *treatment* even when the study isn't an experiment and they're comparing regular populations. Hey, don't blame me! I'm just following the proper statistical terminology.

Here are the general steps in a one-way ANOVA:

- 1. Check the ANOVA conditions, using the data collected from each of the k populations.**

See the next section, "Checking the Conditions" for the specifics on these conditions.

- 2. Set up the hypotheses $H_0: \mu_1 = \mu_2 = \dots = \mu_k$ versus $H_a: \text{At least two of the population means are different}$.**

Another way to state your alternative hypothesis is by saying $H_a: \text{At least two of } \mu_1, \mu_2, \dots, \mu_k \text{ are different}$.

- 3. Collect data from k random samples, one from each population.**

- 4. Conduct an F -test on the data from step three, using the hypotheses from step two, and find the p -value.**

See the section "Doing the F -Test" later in this chapter for these instructions.

- 5. Make your conclusions: If you reject H_0 (when your p -value is less than 0.05 or your predetermined α level), you conclude that at least two of the population means are different; otherwise, you conclude that you didn't have enough evidence to reject H_0 (you can't say the means are different).**

If these steps seem like a foreign language to you, don't fear — I describe each in detail in the sections that follow.

Checking the Conditions

Step one of ANOVA is checking to be sure all necessary conditions are met before diving into the data analysis. The conditions for using ANOVA are just an extension of the conditions for a *t*-test (see the section “Comparing Two Means with a *t*-Test”). The following conditions all need to hold in order to conduct ANOVA:

- ✓ The k populations are independent. In other words, their outcomes don’t affect each other.
- ✓ The k populations each have a normal distribution.
- ✓ The variances of the k normal distributions are equal.

Verifying independence

To check the first condition, examine how the data were collected from each of the separate populations. In order to maintain independence, the outcomes from one population can’t affect the outcomes of the other populations. If the data have been collected by using a separate random sample from each population (*random* here meaning that each individual in the population had an equal chance of being selected), this factor ensures independence at the strongest level.

In the watermelon seed-spitting data (see Table 9-1), the data aren’t randomly sampled from each age group because the data represent everyone who participated in the contest. But, you can argue that in most cases the seed-spitting distances from one age group don’t affect the seed-spitting distances from the other age groups, so the independence assumption is relatively okay.

Looking for what’s normal

The second ANOVA condition is that each of the k populations has a normal distribution. To check this condition, make a separate histogram of the data from each group and see whether it resembles a normal distribution. Data from a normal distribution should look symmetric (in other words, if you split the histogram down the middle, it looks the same on each side) and have a bell shape. Don’t expect the data in each histogram to follow a normal distribution exactly (remember, it’s only a sample), but it shouldn’t be extremely different from a normal, bell-shaped distribution.

Because the seed-spitting data contain only five children per age group, checking conditions can be iffy. But in this case, you have past years’ data for 200 children in each age group, so you can use that to check the conditions. The histograms and descriptive statistics of the seed-spitting data for the four age groups are shown in Figure 9-2, all in one panel, so you can easily compare them to each other on the same scale.

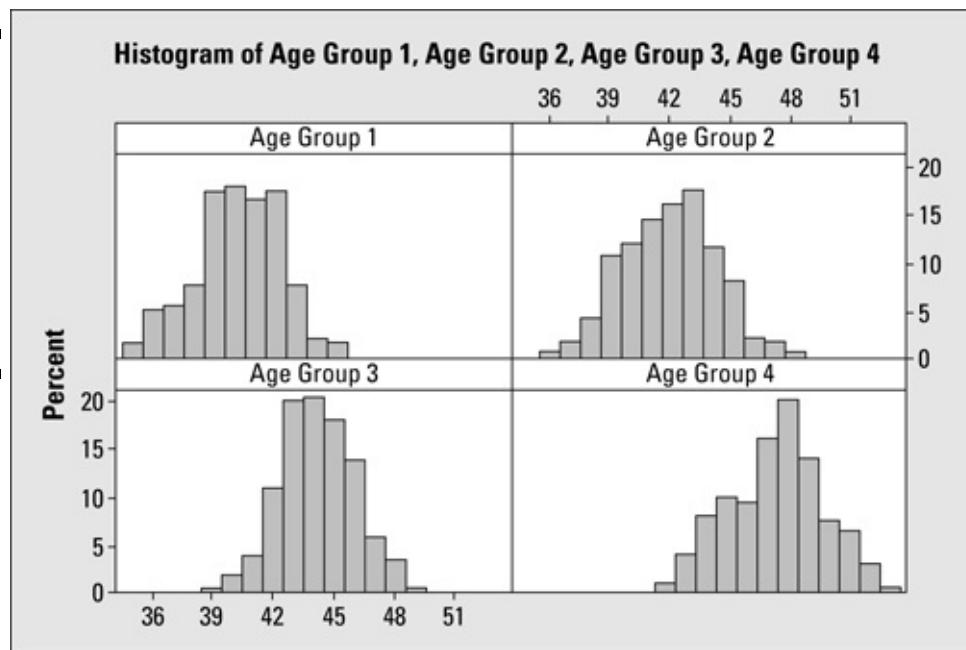
Looking at the four histograms in Figure 9-2, you can see that each graph resembles a bell shape; the normality condition isn’t being severely violated here. (Red flags should come up if you see two peaks in the data, a skewed shape where the peak is off to one side, or a

flat histogram, for example.)



You can use Minitab to make histograms for each of your samples and have all of them appear on one large panel, all using the same scale. To do this, go to Graph>Histogram and click OK. Choose the variables that represent data from each sample by highlighting them in the left-hand box and clicking Select. Then click on Multiple Graphs, and a new window opens. Under the Show Graph Variables option, check the following box: In separate panels of the same graph. On the Same Scales for Graphs option, check the box for x and the box for y . This option gives you the same scale on both the x and y axes for all the histograms. Then click OK.

Figure 9-2:
Checking
ANOVA
conditions by
using
histograms
and
descriptive
statistics.



Descriptive Statistics: Age Group 1, Age Group 2, Age Group 3, Age Group 4

Variable	Total	Count	Mean	Variance
Age Group 1	200	40.116	4.256	
Age Group 2	200	41.880	4.994	
Age Group 3	200	44.165	3.249	
Age Group 4	200	47.405	5.154	

Taking note of spread

The third condition for ANOVA is that the variance in each of the k populations is the same; statisticians call this the *equal variance condition*. You have two ways to check this condition on your data:

- ✓ Calculate each of the variances from each sample and see how they compare.
- ✓ Create one graph showing all the boxplots of each sample sitting side by side. This type of graph is called a *side-by-side boxplot*. (See your Stats I text or my book *Statistics For Dummies* (Wiley) for more information on boxplots.)

If one or more of the calculated variances is significantly different from the others, the equal variance condition is not likely to be met. What does *significantly different* mean? A hypothesis test for equal variances is the statistical tool used to handle this question; however, it falls outside the scope of most Stats II courses, so for now you can make a judgment call. I always say that if the differences in the calculated variances are enough for you to write home about (say they differ by 10 percent or more), the equal variance condition is likely not to be met.

Similarly, if the lengths of one or more of the side-by-side boxplots looks different enough for you to write home about, the equal variance condition is not likely to be met. (But listen, if you really do write home about any of your statistical issues, you may want to spice up your life a bit.)



The length of the box portion of a boxplot is called the *interquartile range*. You calculate it by taking the third quartile (the 75th percentile) minus the first quartile (the 25th percentile.) See your Stats I text or *Statistics For Dummies* for more info.

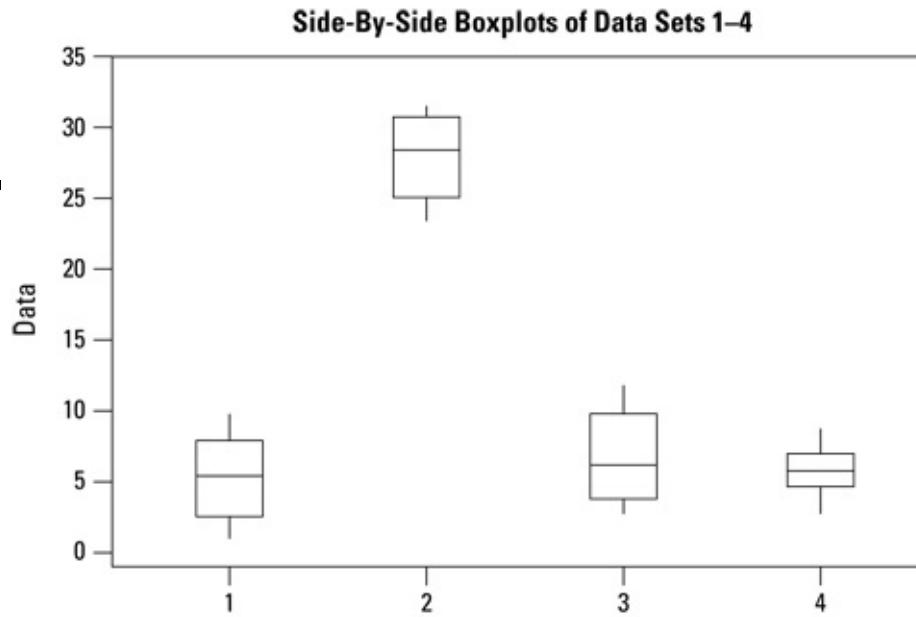
Table 9-2 shows an example of four small data sets with each of their calculated variances shown in the last row. Note that the variance of Data Set 4 is significantly smaller than the others. In this case, it's safe to say that the equal variance condition is not met.

Table 9-2 Comparing Variances of Four Data Sets to Check the Equal Variance Condition

Data Set 1	Data Set 2	Data Set 3	Data Set 4
1	32	4	3
2	24	3	4
3	27	5	5
4	32	10	5
5	31	7	6
6	28	4	6
7	30	8	7
8	26	12	7
9	31	9	8
10	24	10	9
Variance = 9.167	Variance = 9.833	Variance = 9.511	Variance = 3.333

Figure 9-3 shows the side-by-side boxplots for these same four data sets. You see that the boxplot for Data Set 4 has an interquartile range (length of the box) that's significantly smaller than the others. I calculated the actual interquartile ranges for these four data sets; they're 5.50, 5.75, 6.00, and 2.50, respectively. These findings confirm the conclusion that the equal variance condition is not met, due to group 4's much smaller variability.

boxplots to check the equal variance condition.



To find descriptive statistics (including the variance and interquartile range) for each sample, go to Stat>Basic Statistics>Display Descriptive Statistics. Click on each variable in the left-hand box for which you want the descriptive statistics, and click Select. Click on the Statistics option, and a new window appears with tons of different types of statistics. Click on the ones you want and click off the ones you don't want. Click OK. Then click OK again. Your descriptive statistics are calculated.

To find side-by-side boxplots in Minitab, go to Graph>Boxplot. A window appears. Click on the picture for Multiple Y's, Simple, and then click OK. Highlight the variables from the left-hand side that you want to compare, and click Select. Then click OK.



Note that you don't need the sample sizes in each group to be equal to carry out ANOVA; however, in Stats II, you'll typically see what statisticians call a *balanced design*, where each sample from each population has the same sample size. (As I explain in Chapter 3, for more precision in your data, the larger the sample sizes, the better.)

For the seed-splitting data, the variances for each age group are listed in Figure 9-2. These variances are close enough to say the equal variance condition is met.

Setting Up the Hypotheses

Step two of ANOVA is setting up the hypotheses to be tested. You're testing to see if all the population means can be deemed equal to each other. The null hypothesis for ANOVA is that all the population means are equal. That is, $H_0: \mu_1 = \mu_2 = \dots = \mu_k$, where μ_1 is the mean of the first population, μ_2 is the mean of the second population, and so on.

until you reach μ_k (the mean of the k^{th} population).



What appears in the alternative hypothesis (H_a) must be the opposite of what's in the null hypothesis (H_0). What's the opposite of having all k of the population's means equal to each other? You may think the opposite is that they're all different. But that's not the case. In order to blow H_0 wide open, all you need is for at least two of those means to not be equal. So, the alternative hypothesis, H_a , is that at least two of the population means are different from each other. That is, H_a : At least two of $\mu_1, \mu_2, \dots, \mu_k$ are different.

Note that H_0 and H_a for ANOVA are an extension of the hypotheses for a two-sample t -test (which only compares two independent populations). And even though the alternative hypothesis in a t -test may be that one mean is greater than, less than, or not equal to the other, you don't consider any alternative other than \neq in ANOVA. (Statisticians use more in-depth models for the others. Aren't you glad someone else is doing it?)



You only want to know whether or not the means are equal — at this stage of the game anyway. After you reach the conclusion that H_0 is rejected in ANOVA, you can proceed to figure out how the means are different, which ones are bigger than others, and so on, using multiple comparisons. Those details appear in Chapter 10.

Doing the F-Test

Step three of ANOVA is collecting the data, and it includes taking k random samples, one from each population. Step four of ANOVA is doing the F -test on this data, which is the heart of the ANOVA procedure. This test is the actual hypothesis test of $H_0: \mu_1 = \mu_2 = \dots = \mu_k$ versus H_a : At least two of $\mu_1, \mu_2, \dots, \mu_k$ are different.

You have to carry out three major steps in order to complete the F -test. **Note:** Don't get these steps confused with the main ANOVA steps; consider the F -test a few steps within a step:

- 1. Break down the variance of y into sums of squares.**
- 2. Find the mean sums of squares.**
- 3. Put the mean sums of squares together to form the F -statistic.**

I describe each step of the F -test in detail and apply it to the example of comparing watermelon seed-spitting distances (see Table 9-1) in the following sections.



Data analysts rely heavily on computer software to conduct each step of the F -test, and you can do the same. All computer software packages organize and summarize the important information from the F -test into a table format for you.

This table of results for ANOVA is called (what else?) the *ANOVA table*. Because the ANOVA table is a critical part of the entire ANOVA process, I start the following sections out by describing how to run ANOVA in Minitab to get the ANOVA table, and I continue to reference this section as I describe each step of the ANOVA process.

Running ANOVA in Minitab



In using Minitab to run ANOVA, you first have to enter the data from the k samples. You can enter the data in one of two ways:

- ✓ **Stacked data:** You enter all the data into two columns. Column one includes the number indicating what sample the data value is from (1 to k), and the responses (y) are in column two. To analyze this data, go to Stat>ANOVA>One-Way Stacked. Highlight the response (y) variable, and click Select. Highlight the factor (population) variable, and click Select. Click OK.
- ✓ **Unstacked data:** You enter data from each sample into a separate column. To analyze the data entered this way, go to Stat>ANOVA>One-Way Unstacked. Highlight the names of the columns where your data are located, and click OK.

I typically use the unstacked version of data entry just because I think it helps visualize the data. However, the choice is up to you, and the results come out the same no matter which method you choose, as long as you're consistent.

Breaking down the variance into sums of squares

Step one of the F -test is splitting up the variability in the y variable into portions that define where the variability is coming from. Each portion of variability is called a sum of squares. The term *analysis of variance* is a great description for exactly how you conduct a test of k population means. With the overall goal of testing whether k population (or treatment) means are equal, you take a random sample from each of the k populations.

You first put all the data together into one big group and measure how much total variability there is; this variability is called the *sums of squares total*, or SSTO. If the data are really diverse, SSTO is large. If the data are very similar, SSTO is small.

You can split the total variability in the combined data set (SSTO) into two parts:

- ✓ **SST:** The variability between the groups, known as the *sums of squares for treatment*

✓ **SSE:** The variability within the groups, known as the *sums of squares for error*

Splitting up the variability in your data results in one of the most important equalities in ANOVA:

$$\text{SSTO} = \text{SST} + \text{SSE}$$



The formula for SSTO is the numerator of the formula for s^2 , the variance of a

single data set, so $\text{SSTO} = \sum_i \sum_j (x_{ij} - \bar{x})^2$, where i and j represent the j^{th} value in the sample from the i^{th} population and \bar{x} is the *overall sample mean* (the mean of the entire data set). So, in terms of ANOVA, SSTO is the total squared distance between the data values and their overall mean.

The formula for SST is $\sum_i n_i (\bar{x}_i - \bar{x})^2$, where n_i is the size of the sample coming from the i^{th} population and \bar{x} is the overall sample mean. SST represents the total squared distance between the means from each sample and the overall mean.

The formula for SSE is $\sum_i \sum_j (x_{ij} - \bar{x}_i)^2$, where x_{ij} is the j^{th} value in the sample from the i^{th} population and \bar{x}_i is the mean of the sample coming from the i^{th} population. This formula represents the total squared distance between the values in each sample and their corresponding sample means. Using algebra, you can confirm (with some serious elbow grease) that $\text{SSTO} = \text{SST} + \text{SSE}$.



The Minitab output for the watermelon seed-spitting contest for the four age groups is shown in Figure 9-4. Under the Source column of the ANOVA table, you see Factor listed in row one. The factor variable (as described by Minitab) represents the treatment or population variable. In column three of the Factor row, you see the SST, which is equal to 89.75. In the Error row (row two), you locate the SSE in column three, which equals 56.80. In column three of the Total row (row three), you see the SSTO, which is 146.55. Using the values of SST, SSE, and SSTO from the Minitab output, you can verify that $\text{SST} + \text{SSE} = \text{SSTO}$.

Figure 9-4: One-Way ANOVA: Age Group 1, Age Group 2, Age Group 3, Age Group 4

Source	DF	SS	MS	F	P
Factor	3	89.75	29.92	8.43	0.001
Error	16	56.80	3.55		
Total	19	146.55			

$S = 1.884$ $R-\text{Sq} = 61.24\%$ $R-\text{Sq}(\text{adj}) = 53.97\%$

Now you're ready to use these sums of squares to complete the next step of the *F*-test.

Locating those mean sums of squares

After you have the sums of squares for treatment, SST, and the sums of squares for error, SSE (see the preceding section for more on these), you want to compare them to see whether the variability in the y -values due to the model (SST) is large compared to the amount of error left over in the data after the groups have been accounted for (SSE). So you ultimately want a ratio that somehow compares SST to SSE.

To make this ratio form a statistic that they know how to work with (in this case, an F -statistic), statisticians decided to find the means of SST and SSE and work with that. Finding the mean sums of squares is the second step of the F -test, and the mean sums are as follows:

- ✓ **MST** is the *mean sums of squares for treatments*, which measures the mean variability that occurs between the different treatments (the different samples in the data). What you're looking for is the amount of variability in the data as you move from one sample to another. A great deal of variability between samples (treatments) may indicate that the populations are different as well.

You can find MST by taking SST and dividing by $k - 1$ (where k is the number of treatments).

- ✓ **MSE** is the *mean sums of squares for error*, which measures the mean within-treatment variability. The *within-treatment variability* is the amount of variability that you see within each sample itself, due to chance and/or other factors not included in the model.

You can find MSE by taking SSE and dividing by $n - k$ (where n is the total sample size and k is the number of treatments). The values of $k - 1$ and $n - k$ are called the *degrees of freedom* (or df) for SST and SSE, respectively.



Minitab calculates and posts the degrees of freedom for SST, SSE, MST, and MSE in the ANOVA table in columns two and four, respectively.

From the ANOVA table for the seed-spitting data in Figure 9-4, you can see that column two has the heading DF, which stands for degrees of freedom. You can find the degrees of freedom for SST in the Factor row (row two); this value is equal to $k - 1 = 4 - 1 = 3$. The degrees of freedom for SSE is found to be $n - k = 20 - 4 = 16$. (Remember, you have four age groups and five children in each group for a total of $n = 20$ data values.) The degrees of freedom for SSTO is $n - 1 = 20 - 1 = 19$ (found in the Total row under DF). You can verify that the degrees of freedom for SSTO = degrees of freedom for SST + degrees of freedom for SSE.

The values of MST and MSE are shown in column four of Figure 9-4, with the heading MS. You can see the MST in the Factor row, which is 29.92. This value was calculated by

taking SST = 89.75 and dividing it by degrees of freedom, 3. You can see MSE in the Error row, equal to 3.55. MSE is found by taking SSE = 56.80 and dividing it by its degrees of freedom, 16.

By finding the mean sums of squares, you've completed step two of the F-test, but don't stop here! You need to continue to the next section in order to complete the process.

Figuring the F-statistic

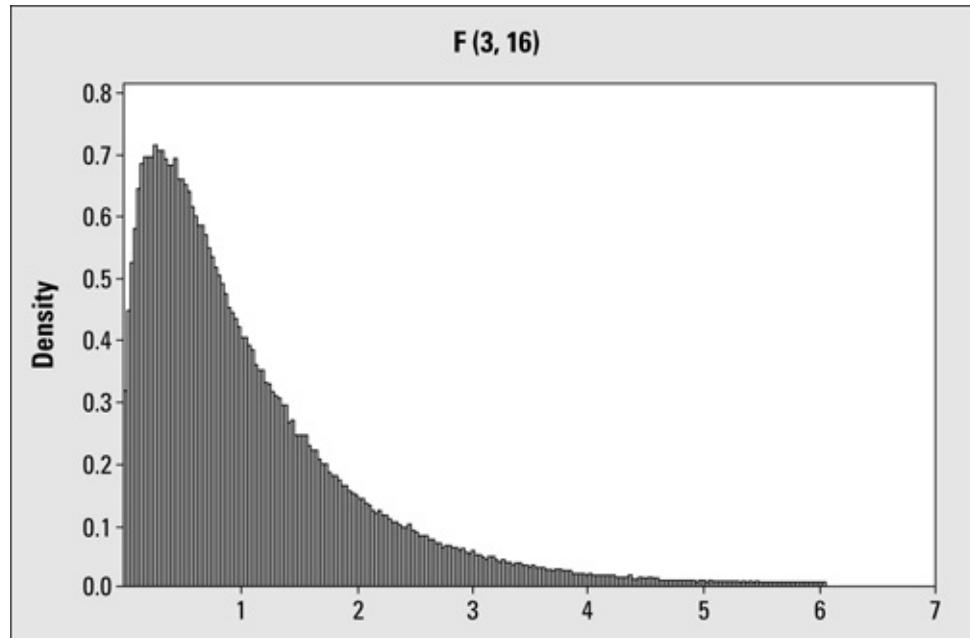
The test statistic for the test of the equality of the k population means is $F = \frac{MST}{MSE}$. The result of this formula is called the *F-statistic*. The *F*-statistic has an *F*-distribution, which is equivalent to the square of a *t*-test (when the numerator degrees of freedom is 1; see more on this interesting connection between the *t*- and *F*-distributions in Chapter 12). All *F*-distributions start at zero and are skewed to the right. The degree of curvature and the height of the curvature of each *F*-distribution is reflected in two degrees of freedom, represented by $k - 1$ and $n - k$. (These come from the denominators of MST and MSE, respectively, where n is the total sample size and k is the total number of treatments or populations.) A shorthand way of denoting the *F*-distribution for this test is $F_{(k-1, n-k)}$.

In the watermelon seed-spitting example, you're comparing four means and have a sample size of five from each population. Figure 9-5 shows the corresponding *F*-distribution, which has degrees of freedom $4 - 1 = 3$ and $20 - 4 = 16$; in other words $F_{(3, 16)}$.



You can see the *F*-statistic on the Minitab ANOVA output (see Figure 9-4) in the Factor row, under the column indicated by F. For the seed-spitting example, the value of the *F*-statistic is 8.43. This number was found by taking MST = 29.92 divided by MSE = 3.55. Then locate 8.43 on the *F*-distribution in Figure 9-5 to see where it stands in terms of its *p*-value. (Turns out it's waaay out there; more on that in the next section.)

Figure 9-5: *F*-distribution with $(3, 16)$ degrees of freedom.



Be sure to not exchange the order of the degrees of freedom for the F -distribution. The difference between $F_{(3, 16)}$ and $F_{(16, 3)}$ is a big one.

Making conclusions from ANOVA

If you've completed the F -test and found your F -statistic (step four in the ANOVA process), you're ready for step five of ANOVA: making conclusions for your hypothesis test of the k population means. If you haven't already done so, you can compare the F -statistic to the corresponding F -distribution with $(k - 1, n - k)$ degrees of freedom to see where it stands and make a conclusion. You can make the conclusion in one of two ways: the p -value approach or the critical-value approach. The approach you use depends primarily on whether you have access to a computer, especially during exams. I describe these two approaches in the following sections.

Using the p -value approach



On Minitab ANOVA output (see Figure 9-4), the value of the F -statistic is located in the Factor row, under the column noted by F. The associated p -value for the F -test is located in the Factor row under the column headed by P. The p -value tells you whether or not you can reject H_0 .

- ✓ **If the p -value is less than your predetermined α (typically 0.05), reject H_0 .**
Conclude that the k population means aren't all equal and that at least two of them are different.
- ✗ **If the p -value is greater than α , then you can't reject H_0 .** You don't have enough evidence in your data to say the k population means have any differences.

The F -statistic for comparing the mean watermelon seed-spitting distances for the four age groups is 8.43. The p -value as indicated in Figure 9-4 is 0.001. That means the results are highly statistically significant. You reject H_0 and conclude that at least one pair of age groups differs in its mean watermelon seed-spitting distances. (You would hope that a 17-year-old could do a lot better than a 6-year-old, but maybe those 6-year-olds have a lot more spitting practice than 17-year-olds do.)

Using Figure 9-5, you see how the F -statistic of 8.43 stands on the F -distribution with $(4 - 1, 20 - 4) = (3, 16)$ degrees of freedom. You can see that it's way off to the right, out of sight. It makes sense that the p -value, which measures the probability of being beyond that F -statistic, is 0.001.

Using critical values

If you're in a situation where you don't have access to a computer (as is still the case in many statistics courses today when it comes to taking exams), finding the exact p -value for the F -statistic isn't possible using a table. You just choose the p -value of the F -statistic that's closest to yours. However, if you do have access to a computer while doing homework or an exam, statistical software packages automatically calculate all p -values exactly, so you can see them on any computer output.



To approximate the p -value from your F -statistic in the event you don't have a computer or computer output available, you find a cutoff value on the F -distribution with $(k - 1, n - k)$ degrees of freedom that draws a line in the sand between rejecting H_0 and not rejecting H_0 . This cutoff, also known as the *critical value*, is determined by your predetermined α (typically 0.05). You choose the critical value so that the area to its right on the F -distribution is equal to α .

F -distribution tables are available in various statistics textbooks and Web sites for other values of α ; however, $\alpha = 0.05$ is by far the most common α level used for the F -distribution and is sufficient for your purposes.

This table of values for the F -distribution is called the *F -table*, and students typically receive these with their exams. For the seed-spitting example, the F -statistic has an F -distribution with degrees of freedom $(3, 16)$, where $3 = k - 1$, and $16 = n - k$. To find the critical value, consult an F -table (Table A-5 in the appendix). Look up the degrees of freedom $(3, 16)$, and you'll find that the critical value is 3.2389 (or 3.24). Your F -statistic for the seed-spitting example is 8.43, which is well beyond this critical value (you can see how 8.43 compares to 3.24 by looking at Figure 9-5). Your conclusion is to reject H_0 at level α . At least two of the age groups differ on mean seed-spitting distances.



With the critical value approach, any F -statistic that lies beyond the critical value results in rejecting H_0 , no matter how far from or close to the line it is. If your F -statistic is beyond the value found in the F -table you consult, then you reject H_0 and say at least two of the treatments (or populations) have different means.

What's next?

After you've rejected H_0 in the F -test and concluded that not all the populations means are the same, your next question may be: Which ones are different? You can answer that question by using a statistical technique called *multiple comparisons*. Statisticians use many different multiple comparison procedures to further explore the means themselves after the F -test has been rejected. I discuss and apply some of the more common multiple comparison techniques in Chapter 10.

Checking the Fit of the ANOVA Model

As with any other model, you must determine how well the ANOVA model fits before you can use its results with confidence. In the case of ANOVA, the model basically boils down to a treatment variable (also known as the population you're in) plus an error term. To assess how well that model fits the data, see the values of R^2 and R^2 adjusted on the last line of the ANOVA output below the ANOVA table. For the seed-spitting data, you see those values at the bottom of Figure 9-4.

- ✓ **The value of R^2 measures the percentage of the variability in the response variable (y) explained by the explanatory variable (x).** In the case of ANOVA, the x variable is the factor due to treatment (where the treatment can represent a population being compared). A high value of R^2 (say, above 80 percent) means this model fits well.
- ✓ **The value of R^2 adjusted, the preferred measure, takes R^2 and adjusts it for the number of variables in the model.** In the case of one-way ANOVA, you have only one variable, the factor due to treatment, so R^2 and R^2 adjusted won't be very far apart. For more on R^2 and R^2 adjusted, see Chapter 6.

For the watermelon seed-spitting data, the value of R^2 adjusted (as found in the last row of Figure 9-4) is only 53.97 percent. That means age group (shown to be statistically significant by the F -test; see the section "Making conclusions from ANOVA") explains just over half of the variability in the watermelon seed-spitting distances. Because of that connection, you may find other variables you can examine in addition to age group, making an even better model for predicting how far those seeds will go.

As you see in Figure 9-1, the results of the t -test done to compare the spitting distances of males and females in the section "Comparing Two Means with a t -Test" show that males

and females were significantly different on mean seed-spitting distances (p -value = 0.039 < 0.05). So I would venture a guess that if you include gender as well as age group, thereby creating what statisticians call a *two-factor ANOVA* (or *two-way ANOVA*), the resulting model would fit the data even better, resulting in higher values of R^2 and R^2 adjusted. (Chapter 11 walks you through the two-way ANOVA.)

Upfront rejection is the best policy for most refusal letters

Many medical and psychological studies use designed experiments to compare the responses of several different treatments, looking for differences. A *designed experiment* is a study in which subjects are randomly assigned to treatments (experimental conditions) and their responses are recorded. The results are used to compare treatments to see which one(s) work best, which ones work equally well, and so on.

Ohio State University researchers conducted one such experiment using ANOVA to determine the most effective way to write a rejection letter. (Is there really a best way to say “no” to someone? Turns out the answer is “yes.”) The experiment tested three traditional principles of writing refusal letters:

- ✓ Using a buffer, which is a neutral or positive sentence that delays the negative information
- ✓ Placing the reason before the refusal
- ✓ Ending the letter on a positive note as a way of reselling the business

Subjects were randomly assigned to treatments, and their responses to the rejection letters were compared (likely on some sort of scale such as 1 = very negative to 7 = very positive with 4 being a neutral response).

You can analyze this scenario by using ANOVA because it compares three treatments (forms of the rejection letters) on some quantitative variable (response to the letter). You can argue that response to the letter isn't a continuous variable, however it has enough possible values that ANOVA isn't unreasonable. The data were also shown to have a bell shape.

The null hypothesis would be H_0 : Mean responses to the three types of rejection letters are equal versus H_a : At least two forms of the rejection letter resulted in different mean responses.

In the end, the researchers did find some significant results; the different ways the rejection letter was written affected the participants differently (so the F-test was rejected). Using multiple comparison procedures (see Chapter 10), you could go in and determine which forms of the rejection letters gave different responses and how the responses differed.

In case you have to write a rejection letter at some point, the researchers recommend the following guidelines:

- ✓ Don't use buffers to begin negative messages.

- ✓ Give a reason for the refusal when it makes the sender's boss look good.
- ✓ Present the negative positively but clearly; offer an alternative or compromise if possible.
- ✓ A positive ending isn't necessary.

Sorting Out the Means with Multiple Comparisons

In This Chapter

- ▶ Knowing when and how to follow up ANOVA with multiple comparisons
 - ▶ Comparing two well-known multiple comparison procedures
 - ▶ Taking additional procedures into consideration
-

Imagine this: You’re comparing the means of not two, but k independent populations, and you find out (using ANOVA; see Chapter 9) that you reject H_0 : All the population means are equal, and you conclude H_a : At least two of the population means are different. Now you gotta know — which of those populations are different? Answering this question requires a follow-up procedure to ANOVA called *multiple comparisons*, which makes sense because you want to compare the multiple means you have to see which ones are different.

In this chapter, you figure out when you need to use a multiple comparison procedure. Two of the most well-known multiple comparison procedures are Fisher’s LSD (least significant difference) and Tukey’s test. They can help you answer that burning question: So some of the means are different, but which ones are different? In this chapter, I also tell you about other comparison procedures that you may encounter or want to try.

Note: For those individuals who come up with new multiple comparison procedures, the procedures are generally named after them. (It’s like having a star named after you but less romantic and a whole lot more work.)

Following Up after ANOVA

The main reason folks use ANOVA to analyze data is to find out whether there are any differences in a group of population means. Your null hypothesis is that there are no differences, and the alternative hypothesis is that there’s at least one difference somewhere between two of the means. (Note it doesn’t say that all the means have to be different.)

If it’s established that at least two of the population means are different, the next natural question is: “Okay, which ones are different?” Although this is a very simple-sounding question, it doesn’t have a simple answer. The concept of means being different can be interpreted in hundreds of ways. Is one larger than all the others? Are three pairs of them different from each other and the rest all the same? Statisticians have worked long and hard to come up with a wide range of choices of procedures to explore and find differences of all types in two or more population means. This family of procedures is

called *multiple comparisons*.

This section starts off with an example in which the ANOVA procedure was used and H_0 was rejected, leading you to the next step: multiple comparisons. You then get an overview of how and why multiple comparison procedures work.

Comparing cellphone minutes: An example

Suppose you want to compare the average number of cellphone minutes used per month for various age groups, where the age groups are defined as the following:

- ✓ Group 1: 19 years old and under
- ✓ Group 2: 20–39 years old
- ✓ Group 3: Adult males 40–59 years old
- ✓ Group 4: Adult females 60 years old and over

You collect data on a random sample of ten people from each group (where no one knows anyone else to keep independence), and you record the number of minutes each person used their cellphone in one month. The first ten lines of a hypothetical data set are shown in Table 10-1.

Table 10-1 Cellphone Minutes Used in One Month

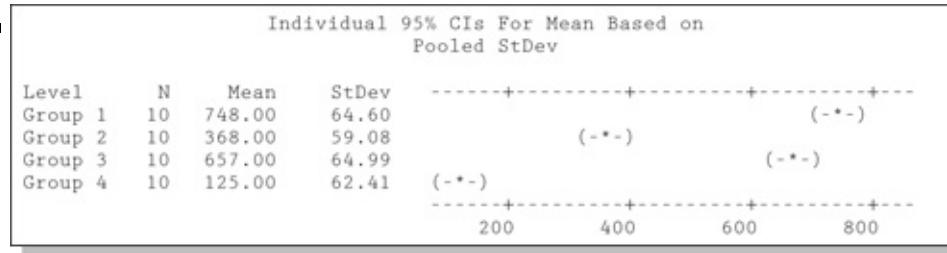
19 and Under (Group 1)	20–39 (Group 2)	40–59 (Group 3)	60 and Over (Group 4)
800	250	700	200
850	350	700	120
800	375	750	150
650	320	650	90
750	430	550	20
680	380	580	150
800	325	700	200
750	410	700	130
690	450	590	160
710	390	650	30

The means and standard deviations of the sample data are shown in Figure 10-1, as well as confidence intervals for each of the population means separately (see Chapter 3 for info on confidence intervals). Looking at Figure 10-1, it appears that all four means are different, with the 19-and-under group heading the pack, 40- to 59-year-olds not far behind, with 20- to 39-year-olds and those over 60 bringing up the rear (in that order).

Knowing that you can't live by sample results alone, you decide that ANOVA is needed to see whether any differences that appear in the samples can be extended to the population (see Chapter 9). By using the ANOVA procedure, you test whether the average cell minutes used is the same across all groups. The results of the ANOVA, using the data

from Table 10-1, are shown in Figure 10-2.

Figure 10-1:
Basic statistics and confidence intervals for the cellphone data.



Looking at Figure 10-2, the F -test for equality of all four population means has a p -value of 0.000, meaning it's less than 0.001. That says at least two of these age groups have a significant difference in their cellphone use (see Chapter 9 for info on the F -test and its results).

Figure 10-2:
ANOVA results for comparing cellphone use for four age groups.

One-way ANOVA: Group 1, Group 2, Group 3, Group 4

Source	DF	SS	MS	F	P
Factor	3	2416010	805337	204.13	0.000
Error	36	142030	3945		
Total	39	2558040			

S = 62.81 R-Sq = 94.5% R-Sq(adj) = 93.99%

Okay, so what's your next question? You just found out that the average number of cellphone minutes used per month isn't the same across these four groups. This doesn't mean all four groups are different (see Chapter 9), but it does mean that at least two groups are significantly different in their cellphone use. So your questions are

- ✓ Which groups are different?
- ✓ How are they different?

Setting the stage for multiple comparison procedures

Determining which populations have differing means after the ANOVA F -test has been rejected involves a new data-analysis technique called *multiple comparisons*. The basic idea of multiple comparison procedures is to compare various means and report where and what the differences are. For example, you may conclude from a multiple comparison procedure that the first population had a mean that was statistically lower than the second population, but it was statistically higher than the mean of the third population.

There are myriad different multiple comparison procedures out there; how do you know which one you should use when? Two basic elements distinguish multiple comparison procedures from each other. I call them purpose and price.

- ✓ **Purpose:** When you know that a group of means aren't all equal, you zoom in to explore the relationships between them, depending on the purpose of your research. Maybe you just want to figure out which means are equivalent and which are not. Maybe you want to sort them into statistically equivalent groups from smallest to largest. Or it may be important to compare the average of one group of

means to the average of another group of means. Different multiple comparison procedures were built for different purposes; for the most part, if you use them for their designed purposes, you have a better chance of finding specific differences you're looking for, if those differences are actually there.

➤ **Price:** Any statistical procedure you use comes with a price: the probability of making a Type I error in your conclusions somewhere during the procedure, due to chance. (A *Type I error* is committed when H_0 is rejected when it shouldn't be; in other words, you think two means are different but they really aren't. See your Stats I textbook or my book *Statistics For Dummies* (Wiley) for more info.) This probability of making at least one Type I error during a multiple comparisons procedure is called the *overall error rate* (also known as the experimentwise error rate (EER), or the familywise error rate). Small overall error rates are of course desirable. Each multiple comparison procedure has its own overall error rate; generally the more specific the relationships are that you're trying to find, the smaller your overall error rate is, assuming you're using a procedure that was designed for your purpose.

In the next section, I describe two all-purpose multiple comparison procedures: Fisher's LSD and Tukey's test.



Don't attempt to explore the data with a multiple comparison procedure if the test for equality of the populations isn't rejected. In this case, you must conclude that you don't have enough evidence to say the population means aren't all equal, so you must stop there. Always look at the p -value of the F -test on the ANOVA output before moving on to conduct any multiple comparisons.

Pinpointing Differing Means with Fisher and Tukey

You've conducted ANOVA to see whether a group of k populations has the same mean, and you rejected H_0 . You conclude that at least two of those populations have different means. But you don't have to stop there; you can go on to find out how many and which means are different by conducting multiple comparison tests.

In this section, you see two of the most well-known multiple comparison procedures: *Fisher's LSD* (also known as *Fisher's protected LSD* or *Fisher's test*) and *Tukey's test* (also known as *Tukey's simultaneous confidence intervals*).



Although I only discuss two procedures in detail in this chapter, tons of other multiple comparison procedures exist (see "So Many Other Procedures, So Little

Time!" at the end of this chapter). Although the other procedures' methods differ a great deal, their overall goal is the same: to figure out which population means differ by comparing their sample means.

Fishing for differences with Fisher's LSD

In this section, I outline the original least significant difference procedure (LSD) and R. A. Fisher's improvement on it (aptly called *Fisher's least significant difference procedure*, or *Fisher's LSD*). The LSD and Fisher's LSD procedures both compare pairs of means using some form of *t*-tests, but they do so in different ways (see Chapter 3 or your Stats I textbook for more on the *t*-test). You also see Fisher's LSD applied to the cellphone example from earlier in this chapter (see the section "Following Up after ANOVA").

The original LSD procedure

To use the original (pre-Fisher) LSD (short for *least significant difference*) simply choose certain pairs of means in advance and conduct a *t*-test on each pair at level $\alpha = 0.05$ to look for differences. LSD doesn't require an ANOVA test first (which is a problem that R. A. Fisher later noticed). If k population means are all to be compared to each other in

pairs using LSD, the number of *t*-tests performed would be represented by $\frac{k(k-1)}{2}$.



Here's how to count the number of *t*-tests when all means are compared. To start, you compare the first mean and the second mean, the first mean and the third mean, and so on until you compare the first mean and the k^{th} mean. Then compare the second and third, second and fourth, and so on all the way down to the $(k - 1)^{\text{th}}$ mean and the k^{th} mean. The total number of pairs of means to compare equals $k * (k - 1)$. Because comparing the two means in either order (mean one and mean two versus mean two and mean one), gives you the same result regarding which one is largest, you divide the total by 2 to avoid double counting. For example, if you have

four populations labeled A, B, C, and D, you have $\frac{4(4-1)}{2} = 6$ *t*-tests to perform: A versus B; A versus C; A versus D; B versus C; B versus D; and C versus D.



The original LSD procedure is very straightforward, easy to conduct, and easy to understand. However, the procedure has some issues. Because each *t*-test is conducted at α level 0.05, each test done has a 5 percent chance of making a Type I error (rejecting H_0 when you shouldn't have, as I explain in Chapter 3).

Although a 5 percent error rate for each test doesn't seem too bad, the errors have a multiplicative effect as the number of tests increases. For example, the chance of making at least one Type I error with six *t*-tests, each at level $\alpha = 0.05$, is 26.50 percent, which is

your overall error rate for the procedure.



If you want or need to know how I arrived at the number 26.50 percent as the overall error rate in that last example, here it goes: The probability of making a Type I error for each test is 0.05. The chance of making at least one error in six tests equals $1 -$ the probability of making no errors in six tests. The chance of not making an error in one test is $1 - \alpha = 0.95$. The chance of no error in six tests is this quantity times itself six times, or $(0.95)^6$, which equals 0.735. Now take $1 -$ this quantity to get $1 - 0.735 = 0.2650$ or 26.50 percent.

Using Fisher's new and improved LSD

R. A. Fisher suggested an improvement over the regular LSD procedure, and his procedure is called *Fisher's LSD*, or *Fisher's protected LSD*. It adds the requirement that an ANOVA *F*-test must be performed first and must be rejected before any pairs of means can be compared individually or collectively. By requiring the *F*-test to be rejected, you're concluding that at least one difference exists in the means. Adding this requirement, the overall error rate of Fisher's LSD is somewhere in the area of α , which is much lower than what you get from the regular LSD procedure.



The downside of Fisher's LSD is that because each *t*-test is made at level α and the overall error rate is also near α , it's good at finding differences that really do exist, but it also makes some false alarms in the process (mainly saying there's a difference when there really isn't).



To conduct Fisher's LSD in Minitab, go to Stat>ANOVA>One-way or One-way unstacked. (If your data appear in two columns with Column 1 representing the population number and Column 2 representing the response, just click One-way because your data are stacked. If your data are shown in k columns, one for each of the k populations, click One-way unstacked.) Highlight the data for the groups you're comparing, and click Select. Then click on Comparisons, and then Fisher's. The individual error rate is listed at 5 (percent), which is typical. If you want to change it, type in the desired error rate (between 0.5 and 0.001), and click OK. You may type in your error rate as a decimal, 0.05, or as a number greater than 1, such as 5. Numbers greater than 1 are interpreted as a percentage.

An ANOVA procedure was done on the cellphone data presented in Table 10-1 to compare the mean number of minutes used for four age groups. Looking at the output in Figure 10-2, you see H_0 (all the population means are equal) was rejected. The next step is to conduct multiple comparisons by using Fisher's LSD to see which population

means differ. Figure 10-3 shows the Minitab output for those tests.

The first block of results shows “Group 1 subtracted from” where Group 1 = age 19 and under. Each line after that represents the other age groups (Group 2 = 20- to 39-year-olds, Group 3 = 40- to 59-year-olds, and Group 4 = age 60 and over). Each line shows the results of comparing the mean for some other group minus the mean for Group 1.

For example, the first row shows Group 2 being compared with Group 1. Moving to the right in that same row, you see the confidence interval for the difference in these two means, which turns out to be -436.97 to -323.03 . Because zero isn’t contained in this interval, you conclude that these two means are different in the populations also. Because this difference ($\mu_2 - \mu_1$) is negative, you also can say that μ_2 is less than μ_1 . Or, a better way to think of it may be that μ_1 is greater than μ_2 . That is, Group 1’s mean is greater than Group 2’s mean.

Figure 10-3:
Output
showing
Fisher’s LSD
applied to the
cellphone
data.

Fisher 95% Individual Confidence Intervals						
All Pairwise Comparisons						
Simultaneous confidence level = 80.32%						
Group 1 subtracted from:						
Group 2	-436.97	-380.00	-323.03	(*-)		
Group 3	-147.97	-91.00	-34.03		(*-)	
Group 4	-679.97	-623.00	-566.03	(*-)		
				-350	0	350 700
Group 2 subtracted from:						
Group 3	232.03	289.00	345.97		(*-)	
Group 4	-299.97	-243.00	-186.03	(-*)		
				-350	0	350 700
Group 3 subtracted from:						
Group 4	-588.97	-532.00	-475.03	(-*)		
				-350	0	350 700



If two means are equal, their difference equals zero, and a confidence interval for the difference should contain zero. If zero isn’t included, you say the means are different.

In this case, each subsequent row in the “Group 1 subtracted from” section of Figure 10-3 shows similar results. None of the confidence intervals contain zero, so you conclude that the mean cellphone use for Group 1 is different from the mean cellphone use for any other group.

Moreover, because all confidence intervals are in negative territory, you can conclude that the mean cellphone use for those users age 19 and under is greater than all the others. (Remember, the mean for this group is subtracted from the others, so a negative difference means its mean is greater).

This process continues as you move down through the output until all six pairs of means are compared to each other. Then you put them all together into one conclusion.

For example, in the second portion of the output, Group 2 is subtracted from Groups 3 and 4. You see the confidence interval for the “Group 3” line is (232.03, 345.97); this gives possible values for Group 3’s mean minus Group 2’s mean. The interval is entirely positive, so conclude that Group 3’s mean is greater than Group 2’s mean (according to this data).

On the next line, the interval for Group 4 minus Group 2 is –299.97 to –186.03. All these numbers are negative, so conclude Group 4’s mean is less than Group 2’s. Combine conclusions to say that Group 3’s mean is greater than Group 2’s, which is greater than Group 4’s.

In the cellphone example, none of the means are equal to each other, and based on the signs of confidence intervals and the results of all the individual pairwise comparisons, the following order of cellphone mean usage prevails: $\mu_1 > \mu_3 > \mu_2 > \mu_4$. (Hypothetical data aside, it may be the case that 40- to 59-year-olds use a lot of cellphone time because of their jobs.) Comparing these results to the sample means in Figure 10-1, this ordering makes sense and the means are separated enough to be declared statistically significant.

Notice near the top of Figure 10-3 that you see “Simultaneous confidence level = 80.32%.” That means the overall error rate for this procedure is $1 - 0.8032 = 0.1968$, which is close to 20 percent, a bit on the high side.

Separating the turkeys with Tukey’s test

The basic idea behind Tukey’s test is to provide a series of simultaneous tests for differences in the means. It still examines all possible pairs of means and keeps the overall error rate at α and also keeps the individual Type I error rate for each pair of means at α . Its distinguishing feature is that it performs the tests all at the same time.

Although the details of the formulas used for Tukey’s test are beyond the scope of this book, they’re not based on the *t*-test but rather something called a *studentized range statistic*, which is based on the highest and lowest means in the group and their difference. The individual error rates are held at 0.05 because Tukey developed a cutoff value for his test statistic that’s based on all pairwise comparisons (no matter how many means are in each group).



To conduct Tukey’s test, go to Stat>ANOVA>One-way or One-way unstacked. (If your data appear in two columns with Column 1 representing the population number and Column 2 representing the response, just click One-way because your data are stacked. If your data are shown in k columns, one for each of the k populations, click One-way unstacked.) Highlight the data for the groups you’re comparing, and click Select. Then click on Comparisons, and then Tukey’s. The familywise (overall) error rate is listed at 5 (percent), which is typical. If you want to change it, type in the desired error rate (between 0.5 and 0.001), and click OK. You

may type in your error rate as a decimal, such as 0.05, or as a number greater than 1, such as 5. Numbers greater than 1 are interpreted as a percentage.

The Minitab output for comparing the groups regarding cellphone use by using Tukey's test appears in Figure 10-4. You can interpret its results in the same ways as those in Figure 10-3. Some of the numbers in the confidence intervals are different, but in this case, the main conclusions are the same: Those age 19 and under use their cellphones most, followed by 40- to 59-year-olds, then 20- to 39-year-olds, and finally those age 60 and over.



The results of Fisher and Tukey don't always agree, usually because the overall error rate of Fisher's procedure is larger than Tukey's (except when only two means are involved). Most statisticians I know prefer Tukey's procedure over Fisher's. That doesn't mean they don't have other procedures they like even better than Tukey's, but Tukey's is a commonly used procedure, and many people like to use it.

Figure 10-4:
Output for
Tukey's test
used to
compare
cellphone
usage.

Tukey 95% Simultaneous Confidence Intervals					
All Pairwise Comparisons					
Individual confidence level = 98.93%					
Group 1 subtracted from:					
Group 2	-455.68	-380.00	-304.32	(- * -)	
Group 3	-166.68	-91.00	-15.32	(- * -)	
Group 4	-698.68	-623.00	-547.32	(- * -)	
				+-----+-----+	
				-700 -350 0 350	
Group 2 subtracted from:					
Group 3	213.32	289.00	364.68	(- * -)	
Group 4	-318.68	-243.00	-167.32	(- * -)	
				+-----+-----+	
				-700 -350 0 350	
Group 3 subtracted from:					
Group 4	-607.68	-532.00	-456.32	(- * -)	
				+-----+-----+	
				-700 -350 0 350	

Examining the Output to Determine the Analysis

Sometimes, the process of answering questions is flipped around in your stats courses. Instead of asking you a question that you use computer output to answer, your professor may give you computer output and ask you to determine the question that the analysis answers. (Kind of like *Jeopardy*.) To work your way backward to the question, you look for clues that tell you what type of analysis was done, and then fill in the details using what you already know about that particular type of analysis.

For example, your professor gives you computer output comparing the ages of ten consumers of each of four cereal brands, labeled C1–C4 (see Figure 10-5). On the analysis, you can see the mean consumer ages for the four cereals being compared to each other, and the analysis also shows and compares the confidence intervals for the averages. The comparison of confidence intervals tells you that you’re dealing with a multiple comparison procedure.

Remember, you’re looking to see whether the confidence intervals for each cereal group overlap; if they don’t, those cereals have different average ages of consumers. If they do overlap, those cereals have mean ages that can’t be declared different.

Based on the data in Figure 10-5, you can see that cereals one (C1) and two (C2) aren’t significantly different, but for cereal three (C3), consumers have a higher average age than C1 and C2. Cereal four (C4) has a significantly higher age than the three others. After the multiple comparison procedure, you know which cereals are different and how they compare to the others.

Figure 10-5:
Multiple
comparison
results for the
cereal
example.

Individual 95% CIs For Mean Based on Pooled StDev				
Level	N	Mean	StDev	
C1	10	8.800	1.687	(---*---)
C2	10	11.800	1.033	(---*---)
C3	10	36.500	7.735	(---*---)
C4	10	55.400	10.309	(---*---)

Sometimes multiple comparison procedures give you groups of means that are equivalent to each other, different from each other, or overlapping. In this case, the final result is $\mu_{\text{C1}} = \mu_{\text{C2}} < \mu_{\text{C3}} < \mu_{\text{C4}}$.

So Many Other Procedures, So Little Time!

Many more multiple comparison procedures exist beyond Fisher’s and Tukey’s imaginations. Those that I discuss in this section are a little more specialized in what they were designed to look for, compared to Tukey’s and Fisher’s. For example, you may want to know whether a certain combination of means is larger than another combination of means; or you may want to only compare specific means to each other, not all the pairs of means.



One thing to note, however, is that in many cases you don’t know exactly what you’re looking for when comparing means — you’re just looking for differences, period. If that’s the case, one of the more general procedures, like Fisher’s or Tukey’s, is the way to go. They’re built for general exploration and do a better job of it than more-specialized procedures.

This section provides an overview of other multiple comparison procedures that exist and tells a little bit about each one, including the people who developed them. Given the dates of when these procedures were developed, I think you'll agree with me that the 1950s was the golden age of the multiple comparison procedure.

Controlling for baloney with the Bonferroni adjustment

The *Bonferroni adjustment* (or *Bonferroni correction*) is a technique used in a host of situations, not just for multiple comparisons. It was basically created to stop people from over-analyzing data. There's a limit to what you should do when analyzing data; there's a line that, when crossed, results in something statisticians call *data snooping*. And the Bonferroni adjustment curbs that.

Data snooping is when someone analyzes her data over and over again until she gets a result that she can say is *statistically significant* (meaning the result is said to have been unlikely to have happened by chance; see Chapter 3). Because the number of tests completed by the data snooper is so high, she's likely to find something significant just by chance. And that result is highly likely to be bogus.

For example, suppose a researcher wants to find out what variable is related to sales of bedroom slippers. He collects data on everything he can think of, including the size of people's feet, the frequency with which they go out to get the paper in their slippers, and their favorite colors. Not finding anything significant, he goes on to examine marital status, age, and income.

Still coming up short, he goes out on a limb and looks at hair color, whether or not the subjects have seen a circus, and where they like to sit on an airplane (aisle or window, sir?). Then wouldn't you know, he strikes gold. Turns out that, according to his data, people who sit on the aisles on planes are more likely to buy bedroom slippers than those who sit by the window or in the middle of a row.

What's wrong with this picture? Too many tests. Each time the researcher examines one variable and conducts a test on it, he chooses an α level at which to conduct the test. (Recall that the α level is the amount of chance you're willing to take of rejecting the null hypothesis and making a false alarm.) As the number of tests increases, the α 's pile up.

Suppose α is chosen to be 0.05. The researcher then has a 5 percent chance of being wrong in finding a significant conclusion, just by chance. So if he does 100 tests, each with a 5 percent chance of an error, on average 5 of those 100 tests will result in a statistically significant result, just by chance. However, researchers who don't know that (or who know and go ahead regardless) find results that they claim are significant even though they're really bogus.



An Italian mathematician named Carlo Emilio Bonferroni (1892–1960) said

“enough already” and created something statisticians call the Bonferroni adjustment in 1950 to control the madness. The *Bonferroni adjustment* simply says that if you’re doing k tests of your data, you can’t do each one at level $\alpha = 0.05$, you need to have an α level for each test equal to $0.05 \div k$.

For example, someone who conducts 20 tests on one data set needs to do each one at level $\alpha = 0.05 \div 20 = 0.0025$. This adjustment makes it harder to find a conclusion that’s significant because the p -value for any test must be less than 0.0025. The Bonferroni adjustment curbs the chance of data snooping until you find something bogus.



The downside of Bonferroni’s adjustment is that it’s very conservative. Although it reduces the chance of concluding two means differ when they really don’t, it fails to catch some differences that really are there. In statistical terms, Bonferroni has power issues. (See your Stats I text or *Statistics For Dummies* for a discussion on power.)

Comparing combinations by using Scheffe’s method

Scheffe’s method was developed in 1953 by Henry Scheffe (1907–1977). This method doesn’t just compare two means at time, like Tukey’s and Fisher’s tests do; it compares all different combinations (called *contrasts*) of the means. For example, if you have the means from four populations, you may want to test to see if their sum equals a certain value, or if the average of two of them equals the average of the two others.

Finding out whodunit with Dunnett’s test

Dunnett’s test was developed in 1955 by Charles Dunnett (1921–1977). *Dunnett’s test* is a special multiple comparison procedure used in a designed experiment that contains a control group. The test compares each treatment group to the control group and determines which treatments do better than others.

Compared to other multiple comparison procedures, Dunnett’s test is better able to find real differences in this situation because it focuses only on the differences between each treatment and the control — not on the differences between every single pair of treatments in the entire study.

Staying cool with Student Newman-Keuls

Student Newman-Keuls test is a different approach from Tukey and Fisher in comparing pairs of means in a multiple comparison procedure. This test comes from the work of three people: “Student,” Newman, and Keuls.

The *Student Newman-Keuls procedure* is based on a stepwise or layer approach. You order

sample means from the smallest to the largest and then examine the differences between the ordered means.

You first test the largest minus smallest difference, and if that turns out to be statistically significant, you conclude that their two respective populations are different in terms of their means. Of the remaining means, the ones that are farthest apart in the order are tested for a significant difference, and so on. You stop when you don't find any more differences.

Duncan's multiple range test

David B. Duncan designed the *Duncan's multiple range test* (MRT) in 1955. The test is based on the Student Newman-Keuls test but has increased power in its ability to detect when the null hypothesis is not true (see Chapter 3) because it increases the value of α at each step of the Student Newman-Keul's test. Duncan's test is used especially in agronomy (crop and farm land management) and other types of agricultural research. One of the neatest things about being a statistician is you never know what kinds of problems you'll be working on or who will use your methods and results.

Although Duncan won the favor of many researchers who used his test (and still do), he wasn't without his critics. Both John Tukey (who developed Tukey's test) and Henry Scheffe (who developed Scheffe's test) accused Duncan's test of being too liberal by not controlling the rate of an overall error (called a *familywise error rate* in the big leagues). But Duncan stood his ground. He said that means are usually never equal anyway, so he wanted to err on the side of making a false alarm (Type I error) rather than missing an opportunity (Type II error) to find out when means are different.



Every procedure in statistics has some chance of making the wrong conclusion, not because of an error in the process but because results vary from data set to data set. You just have to know your situation and choose the procedure that works best for that situation. When in doubt, consult a statistician for help in sorting it all out.

The secret lives of statisticians

Sometimes it's hard to imagine famous people having real lives, and it may be especially hard to picture statisticians doing anything but sitting in the back room calculating numbers. But the truth is, famous statisticians are interesting folks with interesting lives, just like you and me. Consider these stellar statisticians:

✓ **Henry Scheffe:** Scheffe was a very distinguished statistician at University of California, Berkeley. One of his five books *The Analysis of Variance*, written in 1959, is the classic book on the subject and is still used today. (I used it in grad school and still have a copy in my office.) Scheffe enjoyed backpacking, swimming, cycling, reading, and music, having learned to play the recorder during his adult life. Sadly, he died from a bicycle accident on his way to the university in 1977.

✓ **Charles Dunnett:** Nicknamed “Charlie” (did you ever think of famous statisticians as having nicknames?), Dunnett was a distinguished, award-winning professor in the Departments of Mathematics, Statistics, Clinical Epidemiology and Biostatistics at McMaster University in Ontario, Canada. He wrote many papers, two of which were so important that they made it onto the list of the top 25 most-cited statistical papers of all time.

✓ **William Sealy Gosset, or “Student”:** The first name included on the Student Newman-Keuls test is a story in itself. “Student” is a pseudonym of the English statistician William Sealy Gosset (1876–1937). Gosset was a statistician working for the Guinness brewery in Dublin, Ireland, when he became famous for developing the *t*-test, also known as the Student *t*-distribution (see Chapter 3), one of the most commonly used hypothesis tests in the statistical world. Gosset devised the *t*-test as a way to cheaply monitor the quality of beer. He published his work in the best of statistical journals, but his employer regarded his use of statistics in quality control to be a trade secret and wouldn’t let him use his real name on his publications (although all his cronies knew exactly who “Student” was). So if not for Guinness beer, the Student’s *t*-test would have been called the Gossett *t*-test (or you’d be drinking “Gosset beer”).

Going nonparametric with the Kruskal-Wallis test

The Kruskal-Wallis test was developed in 1952 by American statisticians William Kruskal (1919–2005) and W. Allen Wallis (1912–1998). The *Kruskal-Wallis test* is the nonparametric version of a multiple comparison procedure. Nonparametric procedures don’t have nearly as many conditions to meet as their traditional counterparts. All the other procedures described in this chapter require normal distributions from the populations and often the same variance as well.

The Kruskal-Wallis test doesn’t use the actual values of the data; it’s based on ranks (orderings of the data from smallest to largest). The test ranks all the data together, and then looks at how those ranks are distributed out amongst the samples that represent separate populations. If one sample gets all the small ranks, that population is concluded to have a smaller mean than the others, and so on. (Turn to Chapter 16 for the full story on nonparametric statistics and Chapter 19 for all the details of the Kruskal-Wallis test.)

Chapter 11

Finding Your Way through Two-Way ANOVA

In This Chapter

- ▶ Building and carrying out ANOVA with two factors
 - ▶ Getting familiar with (and looking for) interaction effects and main effects
 - ▶ Putting the terms to the test
 - ▶ Demystifying the two-way ANOVA table
-

Analysis of variance (ANOVA) is often used in experiments to see whether different levels of an explanatory variable (x) get different results on some quantitative variable y . The x variable in this case is called a *factor*, and it has certain levels to it, depending on how the experiment is set up.

For example, suppose you want to compare the average change in blood pressure on certain dosages of a drug. The factor is drug dosage. Suppose it has three levels: 10mg per day, 20mg per day, or 30mg per day. Suppose someone else studies the response to that same drug and examines whether the times taken per day (one time or two times) has any effect on blood pressure. In this case, the factor is number of times per day, and it has two levels: once and twice.

Suppose you want to study the effects of dosage *and* number of times taken together because you believe both may have an effect on the response. So what you have is called a *two-way ANOVA*, using two factors together to compare the average response. It's an extension of one-way ANOVA (refer to Chapter 9) with a twist, because the two factors you use may operate on the response differently together than they would separately.

In this chapter, first I give you an example of when you'd need to use a two-way ANOVA. Then I show you how to set up the model, make your way through the ANOVA table, take the F -tests, and draw the appropriate conclusions.

Setting Up the Two-Way ANOVA Model

The two-way ANOVA model extends the ideas of the one-way ANOVA model and adds an interaction term to examine how various combinations of the two factors affect the response. In this section, you see the building blocks of a two-way ANOVA: the treatments, main effects, the interaction term, and the sums of squares equation that puts everything together.

Determining the treatments

The two-way ANOVA model contains two factors, A and B, and each factor has a certain number of levels — say i levels of Factor A and j levels of Factor B.

In the drug study example from the chapter intro, you have A = drug dosage with $i = 1, 2$, or 3, and B = number of times taken per day with $j = 1$ or 2. Each person involved in the study is subject to one of the three different drug dosages and will take the drug in one of the two methods given. That means you have $3 * 2 = 6$ different combinations of Factors A and B that you can apply to the subjects, and you can study these combinations and their effects on blood pressure changes in the two-way ANOVA model.



Each different combination of levels of Factors A and B is called a *treatment* in the model. Table 11-1 shows the six treatments in the drug study. For example, Treatment 4 is the combination of 20mg of the drug taken in two doses of 10mg each per day.

Table 11-1 Six Treatment Combinations for the Drug Study Example

Dosage Amount	One Dose Per Day	Two Doses Per Day
10mg	Treatment 1	Treatment 2
20mg	Treatment 3	Treatment 4
30mg	Treatment 5	Treatment 6



If Factor A has i levels and Factor B has j levels, you have $i * j$ different combinations of treatments in your two-way ANOVA model.

Stepping through the sums of squares

The two-way ANOVA model contains the following three terms:

- ✓ **The main effect A:** Term for the effect of Factor A on the response
- ✓ **The main effect B:** Term for the effect of Factor B on the response
- ✓ **The interaction of A and B:** The effect of the combination of Factors A and B (denoted AB)

The sums of squares equation for the one-way ANOVA (which I cover in Chapter 9) is $SSTO = SST + SSE$, where SSTO is the total variability in the response variable, y ; SST is the variability explained by the treatment variable (call it factor A); and SSE is the variability left over as error.

The purpose of a one-way ANOVA model is to test to see whether the different levels of

Factor A produce different responses in the y variable. The way you do it is by using H_0 : $\mu_1 = \mu_2 = \dots = \mu_i$, where i is the number of levels of Factor A (the treatment variable). If you reject H_0 , Factor A (which separates the data into the groups being compared) is significant. If you can't reject H_0 , you can't conclude that Factor A is significant.



In the two-way ANOVA, you add another factor to the mix (B) plus an interaction term (AB). The sums of squares equation for the two-way ANOVA model is $SSTO = SSA + SSB + SSAB + SSE$. Here, $SSTO$ is the total variability in the y -values; SSA is the sums of squares due to Factor A (representing the variability in the y -values explained by Factor A); and similarly for SSB and Factor B. $SSAB$ is the sums of squares due to the interaction of Factors A and B, and SSE is the amount of variability left unexplained, and deemed error.

Although the mathematical details of all the formulas for these terms are unwieldy and beyond the focus of this book, they just extend the formulas for one-way ANOVA found in Chapter 9. ANOVA handles the calculations for you, so you don't have to worry about that part.



To carry out a two-way ANOVA in Minitab, enter your data in three columns.

- ✓ Column 1 contains the responses (the actual data).
- ✓ Column 2 represents the level of Factor A (Minitab calls it the *row factor*).
- ✓ Column 3 represents the level of Factor B (Minitab calls it the *column factor*).

Go to Stat>Anova>Two-way. Click on Column 1 in the left-hand box, and it appears in the Response box on the right-hand side. Click on Column 2, and it appears in the row factor box; click on Column 3, and it appears in the column factor box. Click OK.

For example, suppose you have six data values in Column 1: 11, 21, 38, 14, 15, and 62. Suppose Column 2 contains 1, 1, 1, 2, 2, 2, and Column 3 contains 1, 2, 3, 1, 2, 3. This means that Factor A has two levels (1, 2), and Factor B has three levels (1, 2, 3). Table 11-2 shows a breakdown of the data values and which combinations of levels and factors are affiliated with them.

Table 11-2 Data and Its Respective Levels from Two Factors

<i>Data Value</i>	<i>Level of Factor A</i>	<i>Level of Factor B</i>
11	1	1
21	1	2
38	1	3
14	2	1

Suppose Factor A has i levels and Factor B has j levels, with a sample of size m collected on each combination of A and B. The degrees of freedom for Factor A, Factor B, and the interaction term AB are $(i - 1)$, $(j - 1)$, and $(i - 1) * (j - 1)$, respectively. This formula is just an extension of the degrees of freedom for the one-way model for Factors A and B. The degrees of freedom for SSTO is $(i * j * m) - 1$, and the degrees of freedom for SSE is $i * j * (m - 1)$. (See Chapter 9 for details on degrees of freedom.)

Understanding Interaction Effects

The interaction effect is the heart of the two-way ANOVA model. Knowing that the two factors may act together in a different way than they would separately is important and must be taken into account. In this section, you see the many ways in which the interaction term AB and the main effects of Factors A and B affect the response variable in a two-way ANOVA model.

What is interaction, anyway?

Interaction is when two factors meet, or interact with each other, on the response in a way that's different from how each factor affects the response separately.

For example, before you can test to see whether dosage of medicine (Factor A) or number of times taken (Factor B) are important in explaining changes in blood pressure, you have to look at how they operate together to affect blood pressure. That is, you have to examine the interaction term.

Suppose you're taking one type of medicine for cholesterol and another medicine for a heart problem. Suppose researchers only looked at the effects of each drug alone, saying each one was good for managing the problem for which it was designed with little or no side effects. Now you come along and mix the two drugs in your system. As far as the individual study results are concerned, all bets are off. With only those separate studies to go on, no one knows how the drugs will interact with each other, and you can find yourself in a great deal of trouble very quickly if you take them together.

Fortunately, drug companies and medical researchers do a great deal of work studying drug interactions, and your pharmacist knows which drugs interact as well. You can bet a statistician was involved in this work from day one!

Baking is another good example of how interaction works. Slurp down one raw egg, drink a cup of milk, and eat a cup of sugar, a cup of flour, and a stick of margarine. Then eat a cup of chocolate chips. Each one of these items has a certain taste, texture, and effect on your taste buds that, in most cases, isn't all that great. But mix them all together in a

bowl and voilà! You have a batch of chocolate chip cookie dough, thanks to the magical effects of interaction.



In any two-way ANOVA, you must check out the interaction term first. If A and B interact with each other and the interaction is statistically significant, you can't examine the effects of either factor separately. Their effects are intertwined and can't be separated.

Interacting with interaction plots

In the two-way ANOVA model, you're dealing with two factors and their interaction. A number of results could come out of this model in terms of significance of the individual terms, as you can see in the following list:

- ✓ Factors A and B are both significant.
- ✓ Factor A is significant but not Factor B.
- ✓ Factor B is significant but not Factor A.
- ✓ Neither Factors A nor B are significant.
- ✓ The interaction term AB is significant, so you don't examine A or B separately.

Figure 11-1 depicts each of these five situations in terms of a diagram using the drug study example. Plots that show how Factors A and B react separately and together on the response variable y are called *interaction plots*. In the following sections, I describe each of these five situations in detail in terms of what the plots tell you and what the results mean in the context of the drug study example.

Factors A and B are significant

Figure 11-1a shows the situation when both A and B are significant in the model and no interaction is present. The lines represent the levels of the times-per-day factor (B); the x -axis represents the levels of the dosage factor (A); and the y -axis represents the average value of the response variable y , which is change in blood pressure, at each combination of treatments.

In order to interpret these interaction plots, you first look at the general trends each line is making. The top line in Figure 11-1a is moving uphill from left to right, meaning that when the drug is taken two times per day, the changes in blood pressure increase as dosage level increases. The bottom line shows a similar result when the drug is taken once per day; blood pressure changes increase as dosage level increases. Assuming these differences are large enough, you conclude that dosage level (Factor A) is significant.

Now you look at how the lines compare to each other. Note that the lines, although parallel, are quite far apart. In particular, the amounts of blood pressure changes are higher overall when taking the drug twice per day (top line) than they are when taking the drug once per day (bottom line). Again, assuming these differences are large enough, you conclude that times per day (Factor B) is significant.

In this case, the different combinations of Factors A and B don't affect the overall trends in blood pressure changes in opposite ways (that is, the lines don't cross each other) so there's no interaction effect between dosage level and times per day.

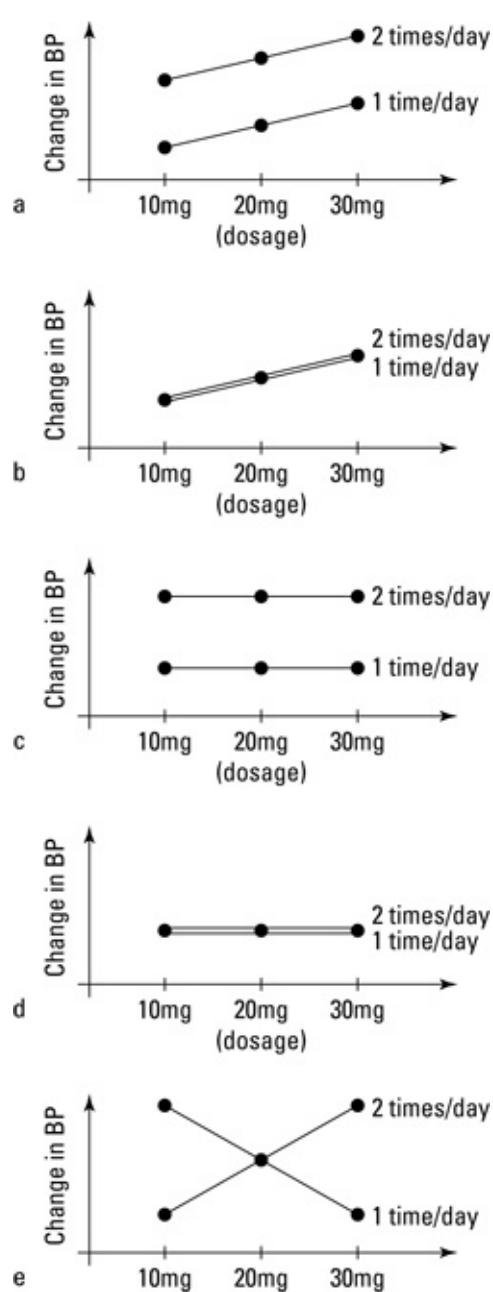


Two parallel lines in an interaction plot indicate a lack of an interaction effect. In other words, the effect of Factor A on the response doesn't change as you move across different levels of Factor B. In the drug study example, the levels of A don't change blood pressure differently for different levels of B.

Factor A is significant but not Factor B

Figure 11-1b shows that blood pressure changes increase across dosage levels for people taking the drug once or twice a day. However, the two lines are so close together that it makes no difference whether you take the drug once or twice a day. So Factor A (dosage) is significant, and Factor B (times per day) isn't. Parallel lines indicate no interaction effect.

Figure 11-1:
Five
examples of
the results
from a two-
way ANOVA
with
interaction.



Factor B is significant but not Factor A

Figure 11-1c shows where Factor B (times per day) is significant but Factor A (dosage level) isn't. The lines are flat across dosage levels, indicating that dosage has no effect on blood pressure. However, the two lines for times per day are spread apart, so their effect on blood pressure is significant. Parallel lines indicate no interaction effect.

Neither factor is significant

Figure 11-1d shows two flat lines that are very close to each other. By the previous discussions about Figures 11-1b and 11-1c, you can guess that this figure represents the case where neither Factor A nor Factor B is significant, and you don't have an interaction effect because the lines are parallel.

Interaction term AB is significant

Finally you get to Figure 11-1e, the most interesting interaction plot of all. The big picture is that because the two lines cross, Factors A and B interact with each other in the way that they operate on the response. If they didn't interact, the lines would be parallel.

Start with the line in Figure 11-1e that increases from left to right (the one for 2 times/day). This line shows that when you take the drug two times per day at the low dose, you get a low change in blood pressure; as you increase dosage, blood pressure change increases also. But when you take the drug once per day, the opposite result happens, as shown by the other line that decreases from left to right in Figure 11-1e.



If you didn't look for a possible interaction effect before you examined the main effects, you may have thought no matter how many times you take this drug per day, the effects will be the same. Not so! Always check out the interaction term first in any two-way ANOVA. If the interaction term is significant, you have no way to pull out the effects due to just factor A or just factor B; they're moot.

Checking the main effects of Factor A or B without checking out the interaction AB term is considered a no-no in the two-way ANOVA world. Another taboo is examining the factors individually (also known as analyzing *main effects*) if the interaction term is significant.

Testing the Terms in Two-Way ANOVA

In a one-way ANOVA, you have only one overall hypothesis test; you use an *F*-test to determine whether the means of the *y* values are the same or different as you go across the levels of the one factor. In two-way ANOVA, you have more items to test besides the overall model. You have the interaction term AB to examine first, and possibly the main effects of A and B. Each test in a two-way ANOVA is an *F*-test based on the ideas of one-way ANOVA (see Chapter 9 for more on this).

To conduct the *F*-tests for these terms, you basically want to see whether more of the total variability in the *y*'s can be explained by the term you're testing compared to what's left in the error term. A large value of *F* means that the term you're testing is significant.

First, you test whether the interaction term AB is significant. To do this, you use the test statistic $F = \frac{MS_{AB}}{MSE}$, which has an *F*-distribution with $(i - 1) * (j - 1)$ degrees of freedom from MS_{AB} (mean sum of squares for the interaction term of A and B) and $i * j * (m - 1)$ degrees of freedom from MSE (mean sum of squares for error), respectively. (Recall that *i* and *j* are the number of levels of A and B, and *m* is the sample size at each combination of A and B.)

If the interaction term isn't significant, you take the AB term out of the model, and you can explore the effects of Factors A and B separately regarding the response variable *y*.

The test for Factor A uses the test statistic $F = \frac{MS_A}{MSE}$, which has an F -distribution with $i - 1$ degrees of freedom from MS_A (mean sum of squares for Factor A) and $i * j * (m - 1)$ degrees of freedom from MSE (mean sum of squares for error), respectively.

Testing for Factor B uses the test statistic $F = \frac{MS_B}{MSE}$, which has an F -distribution with $j - 1$ and $i * j * (m - 1)$ degrees of freedom. (See Chapter 9 for all the details on F-tests, MSE, and degrees of freedom.)



The results you can get from testing the terms of the ANOVA model are the same as those represented in Figure 11-1. They're all provided in Minitab output outlined in the next section, including their sum of squares, degrees of freedom, mean sum of squares, and p -values for their appropriate F -tests.

Running the Two-Way ANOVA Table

The ANOVA table for two-way ANOVA includes the same elements as the ANOVA table for one-way ANOVA (see Chapter 9). But where in the one-way ANOVA you have one line for Factor A's contributions, now you add lines for the effects of Factor B and the interaction term AB. Minitab calculates the ANOVA table for you as part of the output from running a two-way ANOVA.

In this section, you figure out how to interpret the results of a two-way ANOVA, assess the model's fit, and use a multiple comparisons procedure, all using the drug data study.

Interpreting the results: Numbers and graphs

The drug study example involves four people in each treatment combination of three possible dosage levels (10mg, 20mg, and 30mg per day) and two possible times for taking the drug (one time per day and two times per day). The total sample size is $4 * 3 * 2 = 24$. I made up five different data sets in which the analyses represent each of the five scenarios shown in Figure 11-1. Their ANOVA tables, as created by Minitab, are shown in Figure 11-2.

Notice that each ANOVA table in Figure 11-2 shows the degrees of freedom for dosage is $3 - 1 = 2$; the degrees of freedom for times per day is $2 - 1 = 1$; the degrees of freedom for the interaction term is $(3 - 1) * (2 - 1) = 2$; the degrees of freedom for total is $3 * 2 * 4 - 1 = 23$; and the degrees of freedom for error is $3 * 2 * (4 - 1) = 18$.

The order of the graphs in Figure 11-1 and the ANOVA tables in Figure 11-2 isn't the same. Can you match them up? (I promise to give you the answers, so keep reading.)

Here's how the graphs from Figure 11-1 match up with the output in Figure 11-2:

- In the ANOVA table for Figure 11-2a, you see that the interaction term isn't significant (p -value = 0.526), so the main effects can be studied. The p -values for dosage (Factor A) and times taken (Factor B) are 0.000 and 0.001, indicating both Factors A and B are significant; this matches the plot in Figure 11-1a.
- In Figure 11-2b, you see that the p -value for interaction is significant (p -value = 0.000), so you can't examine the main effects of Factors A and B (in other words, don't look at their p -values). This represents the situation in Figure 11-1e.
- Figure 11-2c shows nothing is significant. The p -value for the interaction term is 0.513; p -values for main effects of Factors A (dosage) and B (times taken) are 0.926 and 0.416, respectively. These results coincide with Figure 11-1d.
- Figure 11-2d matches Figure 11-1b. It has no interaction effect (p -value = 0.899); dosage (Factor A) is significant (p -value = 0.000), and times per day (Factor B) isn't (p -value = 0.207).
- Figure 11-2e matches Figure 11-1c. The interaction term, dosage * times per day, isn't significant (p -value = 0.855); times per day is significant with p -value 0.000, but dosage level isn't significant (p -value = 0.855).

Figure 11-2:
ANOVA
tables for the
interaction
plots from
Figure 11-1.

Two-way ANOVA: BP versus Dosage, Times

Source	DF	SS	MS	F	P
Dosage	2	56.3333	28.1667	112.67	0.000
Times	1	4.1667	4.1667	16.67	0.001
Interaction	2	0.3333	0.1667	0.67	0.526
Error	18	4.5000	0.2500		
Total	23	65.3333			

S = 0.5 R-Sq = 93.11% R-Sq(adj) = 91.20%

a

Two-way ANOVA: BP versus Dosage, Times

Source	DF	SS	MS	F	P
Dosage	2	0.0833	0.04167	0.16	0.855
Times	1	0.3750	0.37500	1.42	0.249
Interaction	2	16.7500	8.37500	31.74	0.000
Error	18	4.7500	0.26389		
Total	23	21.9583			

S = 0.5137 R-Sq = 78.37% R-Sq(adj) = 72.36%

b

Two-way ANOVA: BP versus Dosage, Times

Source	DF	SS	MS	F	P
Dosage	2	0.0833	0.04167	0.08	0.926
Times	1	0.3750	0.37500	0.69	0.416
Interaction	2	0.7500	0.37500	0.69	0.513
Error	18	9.7500	0.541667		
Total	23	10.9583			

S = 0.7360 R-Sq = 11.03% R-Sq(adj) = 0.00%

c

Two-way ANOVA: BP versus Dosage, Times

Source	DF	SS	MS	F	P
Dosage	2	36.7500	18.3750	47.25	0.000
Times	1	0.6667	0.6667	1.71	0.207
Interaction	2	0.0833	0.0417	0.11	0.899
Error	18	7.0000	0.3889		
Total	23	44.5000			

S = 7.6236 R-Sq = 84.27% R-Sq(adj) = 79.90%

d

Two-way ANOVA: BP versus Dosage, Times

Source	DF	SS	MS	F	P
Dosage	2	0.0833	0.0417	0.16	0.855
Times	1	12.0417	12.0417	45.63	0.000
Interaction	2	0.0833	0.0417	0.16	0.855
Error	18	4.7500	0.2639		
Total	23	16.9583			

S = 0.5137 R-Sq = 71.99% R-Sq(adj) = 64.21%

e

Assessing the fit

To assess the fit of the two-way ANOVA models, you can use the R^2 adjusted (see Chapter 6). The higher this number is, the better (the maximum is 100 percent or 1.00). Notice that all the ANOVA tables in Figure 11-2 show a fairly high R^2 adjusted except for Figure 11-2c. In this table, none of the terms were significant.

Multiple comparisons

In the case where you find that an interaction effect is statistically significant, you can conduct multiple comparisons to see which combinations of Factors A and B create different results in the response. The same ideas hold here as do for multiple comparisons (covered in Chapter 10), except the tests can be performed on all $i * j$

interactions.



To perform multiple comparisons for a two-way ANOVA by using Minitab, enter your responses (data) in Column 1 (C1), your levels of Factor A in Column 2 (C2), and your levels of factor B in Column 3 (C3). Choose Stat>ANOVA>General Linear Model. In the Responses box, enter your Column 1 variable. In Model, enter C1 <space> C2 <space> C1*C2 (for the main effects and the interaction effect, respectively; here, <space> means leave a space). Click on Comparisons. In Terms, enter Columns 2 and 3. Check the Method you want to use for your multiple comparisons (see Chapter 10), and click OK.

Are Whites Whiter in Hot Water? Two-Way ANOVA Investigates

You use two-way ANOVA when you want to compare the means of n populations that are classified according to two different categorical variables (factors). For example, suppose you want to see how four brands of detergent (Brands A, B, C, D) and water temperature (1 = cold, 2 = warm, 3 = hot) work together to affect the whiteness of dirty t-shirts being washed. (Product-testing groups can use this information as well as the detergent companies to investigate or advertise how a detergent measures up to its competitors.)

Because this question involves two different factors and their effects on some numerical (quantitative) variable, you know that you need to do a two-way ANOVA. You can't assume that water temperature affects whiteness of clothes in the same way for each brand, so you need to include an interaction effect of brand and temperature in the two-way ANOVA model. Because brand of detergent has four possible types (or levels) and water temperature has three possible values (or levels), you have $4 * 3 = 12$ different combinations to examine in terms of how brand and temperature interact. Those combinations are: Brand A in cold water, Brand A in warm water, Brand A in hot water, Brand B in cold water, Brand B in warm water, Brand B in hot water, and so on.

The resulting two-way ANOVA model looks like this: $y = b_i + w_j + bw_{ij} + e$, where b represents the brand of detergent, w represents the water temperature, y represents the whiteness of the clothes after washing, and bw_{ij} represents the interaction of brand i of detergent ($i = A, B, C, D$) and temperature j of the water ($j = 1, 2, 3$). (Note that e represents the amount of variation in the y values (whiteness) that isn't explained by either brand or temperature.)

Suppose you decide to run the experiment five times on each of the 12 combinations, which means 60 observations. (That's 60 t-shirts to wash — hey, it's a dirty job but someone's got to do it!) The results of the two-way ANOVA are shown in Figure 11-3.

Figure 11-3:
ANOVA table
for the
clothing
example.

ANOVA Table: Clothing Example					
Source	DF	SS	MS	F	P
Brand	3	22.983	7.6611	20.89	0.000
Water	2	1.433	0.7167	1.95	0.153
Interaction	6	308.167	51.3611	140.08	0.000
Error	48	17.600	0.3667		
Total	59	350.183			

S = 0.6055 R-Sq = 94.97% R-Sq(adj) = 93.82%

Note that the degrees of freedom (DF) for Brand, Water, Interaction, Error, and Total were arrived at from the following:

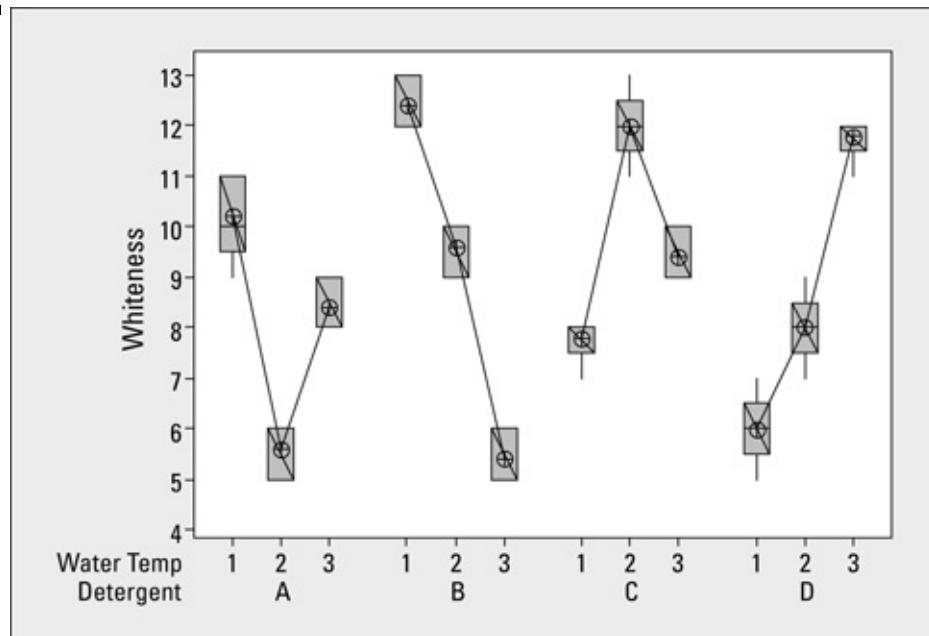
- ✓ DF for brand: $4 - 1 = 3$
- ✓ DF for water temperature: $3 - 1 = 2$
- ✓ DF for interaction term: $(4 - 1) * (3 - 1) = 6$
- ✓ DF for error: $60 - (4 * 3) = 48$
- ✓ DF for total: $n - 1 = 60 - 1 = 59$

Looking at the ANOVA table in Figure 11-3, you can see that the model fits the data very well, with R^2 *adjusted* equal to 93.82 percent. The interaction term (brand of detergent interacting with water temperature) is significant, with a *p*-value of 0.000. This means you can't look separately at the effect of brand of detergent or water temperature separately. One brand of detergent isn't always best, and one water temperature is not always best; it's the combination of the two that has different effects.

Your next question may be: Okay, which combination of detergent brand and water temperature is best? To answer this question, I did multiple comparisons on the means from all 12 combinations. (To do this, I followed the Minitab directions from the previous section.) Luckily, Tukey gives me an overall error rate of only 5 percent, so doing this many tests doesn't lead to making a lot of incorrect conclusions.

Because of the high number of combinations to compare, making sense of all the results on Tukey's output was a little difficult. Instead, I opted to first make boxplots of the data for each combination of brand and water temperature to help me see what was going on. The results of my boxplots are shown in Figure 11-4.

Figure 11-4:
Boxplots
showing how
brand of
detergent
and water
temperature
interact to
affect
clothing
whiteness.



To create one set of boxplots for the data from each of the combinations in a two-way ANOVA, first ask Minitab to conduct a two-way ANOVA (you can find directions in the earlier section “Stepping through the sums of squares”). In that same Minitab window for two-way ANOVA, click Graphs, and a new window comes up. Click Boxplots of Data, and then OK. Finally, click OK to run the analysis and get the boxplots with it.

Figure 11-4 shows four groups of three connected boxes; each group of three represents data from one brand of detergent, tested under each of the three water temperatures (1 = cold, 2 = warm, and 3 = hot). For example, the first group of three shows the data from Brand A under each of the three water temperatures 1, 2, and 3, respectively. Each boxplot shows the results of the whiteness levels for the five shirts washed under that combination of detergent and water temp.

Looking at these plots you can see that each detergent reacts differently with different water temperatures. For example, Brand A does best in cold water (water temp level 1) and worst in warm water (water temp level 2), while Brand C is just the opposite, having the highest scores in warm water and the lowest in cold water. Each detergent does best/worst under a different combination of water temperatures. You can really see why the interaction term in this model is significant!

Now which combination of detergent and water temperature does the best? If you look at the plots, Brand B in cold water looks really good, and so does Brand C in warm water, closely followed perhaps by Brand D in hot water. This is where Tukey’s multiple comparisons come in.

Running multiple comparisons on all 12 combinations of detergent and water temperature, you confirm that the three top combinations identified are all significantly

higher than all the others (because their sample means were higher and their differences from all the other means had p -values less than 0.05). But the top three can't be distinguished from each other (because the p -values for the differences between them all exceed 0.05). Tukey also tells you that the three worst combinations are Brand A in warm water, Brand B in hot water, and Brand D in cold water. And they're all at the bottom of the barrel together (their means are significantly lower than all the rest but can't be distinguished from each other). So no single combination can claim all the bragging rights or shoulder all the blame.

You can imagine the many other comparisons that you could make from here to put the other combinations in some sort of order, but I think the best and worst are the most interesting for this case. It's like the fashion police commenting on what the stars wear on awards night. (Whatever they do wear, let's hope their statistician told them which brand of detergent to use and what water temperature to wash it in!)

Chapter 12

Regression and ANOVA: Surprise Relatives!

In This Chapter

- ▶ Rewriting a regression line as an ANOVA model
- ▶ Connecting regression equations to the ANOVA table

So you're motoring on in your Stats II course, working your way through regression (where you estimate y using one or more x variables; see Chapter 4). Then you hit a new topic, ANOVA, which stands for *analysis of variance* and refers to comparing the means of several populations (see Chapter 9). That seems to be no problem. But wait a minute; now your professor starts talking about how ANOVA is related to regression, and suddenly everything starts to spin out of control. How do you reconcile two techniques that appear to be as different as apples and oranges? That's what this chapter is all about.

Think of this chapter as your bridge across the gap between simple linear regression and ANOVA, allowing you to walk smoothly across, answering any questions that a professor may throw your way. Keep in mind that you don't actually apply these two techniques in this chapter (you can find that information in Chapters 4 and 9); the goal of this chapter is to determine and describe the relationship between regression and ANOVA so they don't look quite so much like an apple and an orange.

Seeing Regression through the Eyes of Variation

Every basic statistical model tries to explain why the different outcomes (y) are what they are. It tries to figure out what factors or explanatory variables (x) can help explain that variability in those y 's. In this section, you start with the y -values by themselves and see how their variability plays a central role in the regression model. This is the first step toward applying ANOVA to the regression model.



No matter what y variable you're interested in predicting, you'll always have variability in those y -values. If you want to predict the length of a fish, for example, you know that fish have many different lengths (indicating a great deal of variability). Even if you put all the fish of the same age and species together, you still have some variability in their lengths (it's less than before but still there nonetheless). The first step in understanding the basic ideas of regression and ANOVA is to understand that variability in the y 's is to be expected, and your job is

to try to figure out what can explain most of it.

Spotting variability and finding an “x-planation”

Both regression and ANOVA work to get a handle on explaining the variability in the y variable using an x variable. After you collect your data, you can find the standard deviation in the y variable to get a sense of how much the data varies within the sample. From there, you collect data on an x variable and see how much it contributes to explaining that variability.

Suppose you notice that people spend different amounts of time on the Internet, and you want to explore why that may be. You start by taking a small sample of 20 people and record how many hours per month they spend on the Internet. The results (in hours) are 20, 20, 22, 39, 40, 19, 20, 32, 33, 29, 24, 26, 30, 46, 37, 26, 45, 15, 24, and 31. The first thing you notice about this data is the large amount of variability in it. The *standard deviation* (average distance from the data values to their mean) of this data set is 8.93 hours, which is quite large given the size of the numbers in the data set.

So you figured out that the y -values — the amount of time someone uses the Internet — have a great deal of variability in them. What can help explain this? Part of the variability is due to chance. But you suspect some variable is out there (call it x) that has some connection to the y variable, and that x variable can help you make more sense out of this seemingly wide range of y -values.

Suppose you have a brainstorm that number of years of education could possibly be related to Internet use. In this case, the explanatory variable (input variable, x) is years of education, and you want to use it to try to estimate y , the number of hours spent on the Internet in a month. You ask a larger random sample of 250 Internet users how many years of education they have (so $n = 250$). You can check out the first ten observations from your data set containing the (x, y) pairs in Table 12-1. If a significant connection of some sort exists between the x -values and the y -values, then you can say that x is helping to explain some of the variability in the y 's. If it explains enough variability, you can place x into a simple regression model and use it to estimate y .

Table 12-1 First Ten Observations from the Education and Internet Use Example

<i>Years of Education</i>	<i>Hours Spent on Internet (In One Month)</i>
15	41
15	32
11	33
10	42
10	28
10	21
10	17

Getting results with regression

After you have a possible x variable picked, you collect pairs of data (x, y) on a random sample of individuals from the population, and you look for a possible linear relationship between them. Looking at the small snippet of 10 out of the 250-person data set in Table 12-1, you can begin to see that you may have a pattern between education and Internet use. It looks like as education increases so does Internet use.

To delve deeper, you make a scatterplot of the data and calculate the correlation (r). If the data appear to follow a straight line (as shown on the scatterplot), go ahead and perform a simple linear regression of the response variable y based on the x variable. The p -value of the x variable in the simple linear regression analysis tells you whether or not the x variable does a significant job in predicting y . (For the details on simple linear regression, see Chapter 4.)



To do a simple linear regression using Minitab, enter your data in two columns: the first column for your x variable and the second column for your y variable (as in Table 12-1). Go to Stat>Regression>Regression. Click on your y variable in the left-hand box; the y variable then appears in the Response box on the right-hand side. Click on your x variable in the left-hand box; the x variable then appears in the Predictor box in the right-hand side. Click OK, and your regression analysis is done. As part of every regression analysis, Minitab also provides you with the corresponding ANOVA results, found at the bottom of the output.

The simple linear regression output that Minitab gives you for the education and Internet example is in Figure 12-1. (Notice the ANOVA output at the bottom; you can see the connection in the upcoming section “Regression and ANOVA: A Meeting of the Models.”)

Regression Analysis: Internet versus Education

The regression equation is
Internet = -8.29 + 3.15 Education

Predictor	Coef	SE Coef	T	P
Constant	-8.290	2.665	-3.11	0.002
Education	3.1460	0.2387	13.18	0.000

S = 7.23134 R-Sq = 41.2% R-Sq(adj) = 41.0%

Analysis of Variance

Source	DF	SS	MS	F	P
Regression	1	9085.6	9085.6	173.75	0.000
Residual Error	248	12968.5	52.3		
Total	249	22054.0			

Figure 12-1:
Output for
simple linear
regression
applied to
education
and Internet
use data.

Looking at Figure 12-1, you see that the p -value on the row marked *Education* is 0.000, which means the p -value's less than 0.001. Therefore the relationship between years of education and Internet use is statistically significant. A scatterplot of the data (not shown here) also indicates that the data appear to have a positive linear relationship, so as you increase number of years of education, Internet use also tends to increase (on average).

Assessing the fit of the regression model

Before you go ahead and use a regression model to make predictions for y based on an x variable, you must first assess the fit of your model. You can do this with a scatterplot and correlation or R^2 .

Using a scatterplot and correlation

One way to get a rough idea of how well your regression model fits is by using a *scatterplot*, which is a graph showing all the pairs of data plotted in the x - y plane. Use the scatterplot to see whether the data appear to fall in the pattern of a line. If the data appear to follow a straight-line pattern (or even something close to that — anything but a curve or a scattering of points that has no pattern at all), you calculate the correlation, r , to see how strong the linear relationship between x and y is. The closer r is to +1 or -1, the stronger the relationship; the closer r is to zero, the weaker the relationship. Minitab can do scatterplots and correlations for you; see Chapter 4 for more on simple linear regression, including making a scatterplot and finding the value of r .



If the data don't have a significant correlation and/or the scatterplot doesn't look linear, stop the analysis; you can't go further to find a line that fits a relationship that doesn't exist.

Using R^2

The more general way of assessing not only the fit of a simple linear regression model but many other models too is to use R^2 , also known as the *coefficient of determination*. (For example, you can use this method in multiple, nonlinear, and logistic regression models in Chapters 5, 7, and 8, to name a few.) In simple linear regression, the value of R^2 (as indicated by Minitab and statisticians as a capital R squared) is equal to the square of the Pearson correlation coefficient, r (indicated by Minitab and statisticians by a small r). In all other situations, R^2 provides a more general measure of model fit. (Note that r only measures the fit of a straight-line relationship between one x variable and one y variable; see Chapter 4.) An even better statistic, R^2 adjusted, modifies R^2 to account for the number of variables in the model. (For more information on R^2 and its use and interpretation, see Chapter 6.)

The value of R^2 adjusted for the model of using education to estimate Internet use (see

Figure 12-1) is equal to 41 percent. This value reflects the percentage of variability in Internet use that can be explained by a person's years of education. This number isn't close to one, but note that r , the square root of 41 percent, is 0.64, which in the case of linear regression indicates a moderate relationship.

This evidence gives you the green light to use the results of the regression analysis to estimate number of hours of Internet use in a month by using years of education. The regression equation as it appears in the top part of the Figure 12-1 output is Internet use = $-8.29 + 3.15 * \text{years of education}$. So if you have 16 years of education, for example, your estimated Internet use is $-8.29 + 3.15 * 16 = 42.11$, or about 42 hours per month (about 10.5 hours per week).

But wait! Look again at Figure 12-1 and zoom in on the bottom part. I didn't ask for anything special to get this info on the Minitab output, but you can see an ANOVA table there. That seems like a fish out of water doesn't it? The next section connects the two, showing you how an ANOVA table can describe regression results (albeit it in a different way).

Regression and ANOVA: A Meeting of the Models

After you've broken down the regression output into all its pieces and parts, the next step toward understanding the connection between regression and ANOVA is to apply the sums of squares from ANOVA to regression (something that's typically not done in a regression analysis). Before you start, think of this process as going to a 3-D movie, where you have to wear special glasses in order to see all the special effects!

In this section, you see the sums of squares in ANOVA applied to regression and how the degrees of freedom work out. You build an ANOVA table for regression and discover how the t -test for a regression coefficient is related to the F -test in ANOVA.

Comparing sums of squares

Sums of squares is a term you may remember from ANOVA (see Chapter 9), but it certainly isn't a term you normally use when talking about regression (see Chapter 4). Yet, you can break down both types of models into sums of squares, and that similarity gets at the true connection between ANOVA and regression.



In step-by-step terms, you first partition out the variability in the y variable by using formulas for sums of squares from ANOVA (sums of squares for total, treatment, and error). Then you find those same sums of squares for regression — this is the twist on the process. You compare the two procedures through their

sums of squares. This section explains how this comparison is done.

Partitioning variability by using SSTO, SSE, and SST for ANOVA

ANOVA is all about partitioning the total variability in the y -values into sums of squares (find all the info you ever need on one-way ANOVA in Chapter 9). The key idea is that $SSTO = SST + SSE$, where $SSTO$ is the total variability in the y -values; SST measures the variability explained by the model (also known as the treatment, or x variable in this case); and SSE measures the variability due to error (what's left over after the model is fit).

Following are the corresponding formulas for $SSTO$, SSE , and SST , where \bar{y} is the mean of the y 's, y_i is each observed value of y , and \hat{y}_i is each predicted value of y from the ANOVA model:

$$SSTO = \sum(y_i - \bar{y})^2$$

$$SSE = \sum(y_i - \hat{y}_i)^2$$

$$SST = \sum(\hat{y}_i - \bar{y})^2$$

Use these formulas to calculate the sums of squares for ANOVA. (Minitab does this for you when it performs ANOVA.) Keep these values of $SSTO$, SST , and SSE . You'll use them to compare to the results from regression.

Finding sums of squares for regression

In regression, you measure the deviations in the y -values by taking each y_i minus its mean, \bar{y} . Square each result and add them all up, and you have $SSTO$. Next, take the residuals, which represent the difference between each y_i and its estimated value from the model, \hat{y}_i . Square the residuals and add them up, and you get the formula for SSE .

After you calculate $SSTO$ and SSE , you need the bridge between them — that is, you need a formula that connects the variability in the y 's ($SSTO$) and the variability in the residuals after fitting the regression line (SSE). That bridge is called the sum of squares for regression, or SSR (equivalent to SST in ANOVA). In regression, \hat{y}_i represents the predicted value of y_i based on the regression model. These are the values on the regression line. To assess how much this regression line helps to predict the y -values, you compare it to the model you'd get without any x variable in it.

Without any other information, the only thing you can do to predict y is look at the average, \bar{y} . So, SSR compares the predicted value from the regression line to the predicted value from the flat line (the mean of the y 's) by subtracting them. The result is $(\hat{y}_i - \bar{y})$. Square each result and sum them all up, and you get the formula for SSR , which is the same as the formula for SST in ANOVA. Voilà!



Instead of calling the sum of squares for the regression model SST as is done in ANOVA, statisticians call it SSR for *sum of squares regression*. Consider SSR to be equivalent to the SST from ANOVA. You need to know the difference because computer output lists the sums of squares for the regression model as SSR, not SST.

To summarize the sums of squares as they apply to regression, you have $SSTO = SSR + SSE$ where

- ✓ SSTO measures the variability in the observed y -values around their mean. This value represents the variance of the y -values.
- ✓ SSE represents the variability between the predicted values for y (the values on the line) and the observed y -values. SSE represents the variability left over after the line has been fit to the data.
- ✓ SSR measures the variability in the predicted values for y (the values on the line) from the mean of y . SSR is the sum of squares due to the regression model (the line) itself.

Minitab calculates all the sums of squares for you as part of the regression analysis. You can see this calculation in the section “Bringing regression to the ANOVA table.”

Dividing up the degrees of freedom

In ANOVA, you test a model for the treatment (population) means by using an F -test, which is $F = \frac{MST}{MSE}$. To get MST (the mean sum of squares for treatment), you take SST (the sum of squares for treatment) and divide by its degrees of freedom. You do the same with MSE (that is, take SSE, the sum of squares for error, and divide by its degrees of freedom). The questions now are, what do those degrees of freedom represent, and how do they relate to regression?

Degrees of freedom in ANOVA

In ANOVA, the degrees of freedom for SSTO is $n - 1$, which represents the sample size minus one. In the formula for SSTO, $\sum(y_i - \bar{y})^2$, you see there are n observed y -values minus one mean. In a very general way, that's where the $n - 1$ comes from.



Note that if you divide SSTO by $n - 1$, you get $\frac{\sum(y_i - \bar{y})^2}{n-1}$, the variance in the y -values. This calculation makes good sense because the variance measures the total variability in the y -values.

Degrees of freedom in regression

The degrees of freedom for SST in ANOVA equal the number of treatments minus one. How does the degrees of freedom idea relate to regression? The number of treatments in regression is equivalent to the number of parameters in a model (a *parameter* being an unknown constant in the model that you're trying to estimate).

When you test a model, you're always comparing it to a different (simpler) model to see whether it fits the data better. In linear regression, you compare your regression line $y = a + bx$, to the horizontal line $y = \bar{y}$. This second, simpler model just uses the mean of y to predict y all the time, no matter what x is. In the regression line, you have two coefficients: one to estimate the parameter for the y -intercept (a) and one to estimate the parameter for slope (b) in the model. In the second, simpler model, you have only one parameter: the value of the mean. The degrees of freedom for SSR in simple linear regression is the difference in the number of parameters from the two models: $2 - 1 = 1$.

The degrees of freedom for SSE in ANOVA is $n - k$. In the formula for SSE, $\sum(\hat{y}_i - \bar{y})^2$, you see there are n predicted y -values, and k is the number of treatments in the model. In regression, the number of parameters in the model is $k = 2$ (the slope and the y -intercept). So you have degrees of freedom $n - 2$ associated with SSE when you're doing regression.



Putting all this together, the degrees of freedom for regression must add up for the equation $SSTO = SSR + SSE$. The degrees of freedom corresponding to this equation are $(n - 1) = (2 - 1) + (n - 2)$, which is true if you do the math. So the degrees of freedom for regression, using the ANOVA approach, all check out. Whew!

In Figure 12-1, you can see the degrees of freedom for each sum of squares listed under the DF column of the ANOVA part of the output. You see SSR has $2 - 1 = 1$ degree of freedom, SSE has $250 - 2 = 248$ degrees of freedom (because $n = 250$ observations were in the data set and $k = 2$ and you find $n - k$ to get degrees of freedom for SSE). The degrees of freedom for SSTO is $250 - 1 = 249$.

Bringing regression to the ANOVA table

In ANOVA, you test your model H_0 : All k population means are equal versus H_a : At least two population means are different by using a F -test. You build your F -test statistic by relating the sums of squares for treatment to the sum of squares for error. To do this, you divide SSE and SST by their degrees of freedom ($n - k$ and $k - 1$, respectively, where n is the sample size and k is the number of treatments) to get the mean sums of squares for error (MSE) and mean sums of squares for treatment (MST). In general, you want MST to be large compared to MSE, indicating that the model fits well. The results of all these statistical gymnastics are summarized by Minitab in a table called (cleverly) the ANOVA table.

The ANOVA table shown in the bottom part of Figure 12-1 for the Internet use data example represents the ANOVA table you get from using the regression line as your model. Under the Source column, you may be used to seeing treatment, error, and total. For regression, the treatment is the regression line, so you see *regression* instead of treatment. The error term in ANOVA is labeled *residual error*, because in regression, you measure error in terms of residuals. Finally you see *total*, which is the same the world around.

The SS column represents the sums of squares for the regression model. The three sums of squares listed in the SS column are SSR (for regression), SSE (for residuals), and SST (total). These sums of squares are calculated using the formulas from the previous section; the degrees of freedom, DF in the table, are found by using the formulas from the previous section also.

The MS column takes the value of SS[you fill in the blank] and divides it by the respective degrees of freedom, just like ANOVA. For example in Figure 12-1, SSE is 12968.5, and the degrees of freedom is 248. Take the first value divided by the second one to get 52.29 or 52.3, which is listed in the ANOVA table for MSE.

The value of the *F*-statistic, using the ANOVA method, is $F = \frac{MST}{MSE} = \frac{9085.6}{52.3} = 173.7$ in the Internet use example, which you can see in column five of the ANOVA part of Figure 12-1 (subject to rounding). The *F*-statistics's *p*-value is calculated based on an *F*-distribution with $k - 1 = 2 - 1 = 1$ and $n - k = 250 - 2 = 248$ degrees of freedom, respectively. (In the Internet use example, the *p*-value listed in the last column of the ANOVA table is 0.000, meaning the regression model fits.) But remember, in regression you don't use an *F*-statistic and an *F*-test. You use a *t*-statistic and a *t*-test. (Whoa . . .)

Relating the *F*- and *t*-statistics: The final frontier

In regression, one way of testing whether the best-fitting line is statistically significant is to test H_0 : Slope = 0 versus H_a : Slope $\neq 0$. To do this, you use a *t*-test (see Chapter 3). The slope is the heart and soul of the regression line, because it describes the main part of the relationship between x and y . If the slope of the line equals zero (you can't reject H_0), you're just left with $y = a$, a horizontal line, and your model $y = a + bx$ isn't doing anything for you.

In ANOVA, you test to see whether the model fits by testing H_0 : The means of the populations are all equal versus H_a : At least two of the population means aren't equal. To do this you use an *F*-test (taking MST and dividing it by MSE; see Chapter 9).



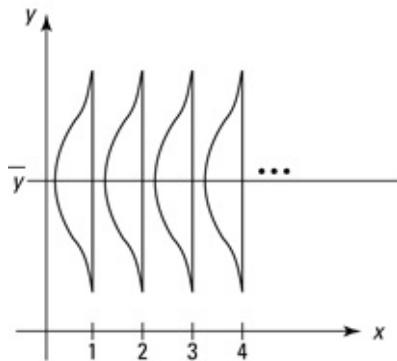
The sets of hypotheses in regression and ANOVA seem totally different, but in essence, they're both doing the same general thing: testing whether a certain model fits. In the regression case, the model you want to see fit is the straight line, and in

the ANOVA case, the model of interest is a set of (normally distributed) populations with at least two different means (and the same variance). Here each population is labeled as a treatment by ANOVA.

But more than that, you can think of it this way: Suppose you took all the populations from the ANOVA and lined them up side by side on an x - y plane (see Figure 12-2). If the means of those distributions are all connected by a flat line (representing the mean of the y 's), then you have no evidence against H_0 in the F -test, so you can't reject it — your model isn't doing anything for you (simply put, it doesn't fit). This idea is similar to the idea of fitting a flat horizontal line through the y -values in regression; a straight-line model with a nonzero slope. This also indicates no relationship between x and y .

The big thing is that statisticians can prove (so you don't have to) that an F -statistic is equivalent to the square of a t -statistic and that the F -distribution is equivalent to the square of a t -distribution when the SSR has $df = 2 - 1 = 1$. And when you have a simple linear regression model, the degrees of freedom is exactly 1! (Note that F is always greater than or equal to zero, which is needed if you're making it the square of something.) So there you have it! The t -statistic for testing the regression model is equivalent to an F -statistic for ANOVA when the ANOVA table is formed for the simple regression model.

Figure 12-2:
Connecting
means of
populations
to the slope
of a line.

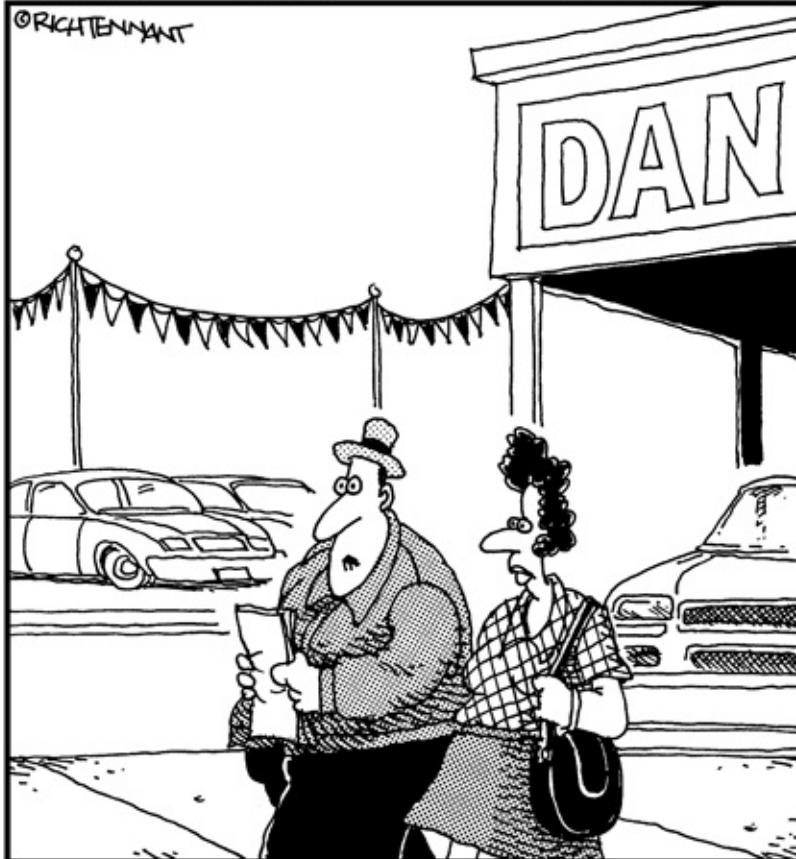


Indeed (the stats professor's way of saying "and this is the *really cool part . . .*"), if you look at the value of the t -statistic for testing the slope of the education variable in Figure 12-1, you see that it's 13.18 (look at the row marked Education and the column marked T). Square that value, and you get 173.71, the F -statistic in the ANOVA table of Figure 12-1. The F -statistic from ANOVA and the square of the t -statistic from regression are equal to each other in Figure 12-2, subject to a little round-off error done by Minitab on the output. (Just like magic! I still get chills just thinking about it.)

Part IV

Building Strong Connections with Chi-Square Tests

The 5th Wave By Rich Tennant



"Is it just me or did the whole '50% satisfaction' statistic seem a little unimpressive?"

In this part . . .

Have you ever wondered if the percentage of M&M'S of each color is the same in every bag? Or whether someone's vote in an election is related to gender? Have you ever wondered if banks really have a case for denying loans based on a low credit score? This part answers all those questions and more using the Chi-square distribution. In particular, you get to use Chi-square to test for independence and to run goodness-of-fit tests.

Chapter 13

Forming Associations with Two-Way Tables

In This Chapter

- ▶ Reading and interpreting two-way tables
 - ▶ Figuring probabilities and checking for independence
 - ▶ Watching out for Simpson’s Paradox
-

Looking for relationships between two categorical variables is a very common goal for researchers. For example, many medical studies center on how some characteristic about a person either raises or lowers his chance of getting some disease. Marketers ask questions like, “Who’s more likely to buy our product: males or females?” Sports stats freaks wonder about things like, “Does winning the coin toss at the beginning of a football game increase a team’s chance of winning the game?” (I believe it does!)

To answer each of the preceding questions, you must first collect data (from a random sample) on the two categorical variables being compared — call them x and y . Then you organize that data into a table that contains columns and rows, showing how many individuals from the sample appear in each combination of x and y . Finally, you use the information in the table to conduct a hypothesis test (called the *Chi-square test*). Using the Chi-square test, you can determine whether you can see a relationship between x and y in the population from which the data were drawn. You need the machinery from Chapter 14 to accomplish this last step.

The goals of this chapter are to help you to understand what it means for two categorical variables (x and y) to be associated and to discover how to use percentages to determine whether a sample data set appears to show a relationship between x and y .

Suppose you’re collecting data on cellphone users, and you want to find out whether more females use cellphones for personal use than males. A study of 508 randomly selected male cellphone users and 508 randomly selected female cellphone users conducted by a wireless company found that women tend to use their phones for personal calls more than men (big shocker). The survey showed that 427 of the women said they used their wireless phones primarily to talk with friends and family, while only 325 of the men admitted to doing so.

But you can’t stop there. You need to break down this information, calculate some percentages, and compare those percentages to see how close they really are. Sample results vary from sample to sample, and differences can appear by chance.

In this chapter, you find out how to organize data from categorical variables (that’s data based on categories rather than measurements) into a table format. This skill is

especially useful when you’re looking for relationships between two categorical variables, such as using a cellphone for personal calls (a yes or no category) and gender (male or female). You also summarize the data to answer your questions. And, finally, you get to figure out, once and for all, what’s going on with that Simpson’s Paradox thing.

Breaking Down a Two-Way Table

A *two-way table* is a table that contains rows and columns that helps you organize data from categorical variables in the following ways:

- ✓ **Rows** represent the possible categories for one categorical variable, such as males and females.
- ✓ **Columns** represent the possible categories for a second categorical variable, such as using your cellphone for personal calls, or not.

Organizing data into a two-way table

To organize your data into a two-way table, first set up the rows and columns. Table 13-1 shows the setup for the cellphone data example that I set up in the chapter introduction.

Table 13-1 Two-Way Table for the Cellphone Data

	<i>Personal Calls: Yes</i>	<i>Personal Calls: No</i>
<i>Males</i>		
<i>Females</i>		

Notice that Table 13-1 has four empty cells inside of it (not counting the empty space in the upper-left corner). Because gender has two choices (male or female) and personal cellphone use has two choices (yes or no), the resulting two-way table has $2 * 2 = 4$ cells.



To figure out the number of cells in any two-way table, multiply the number of possible categories for the row variable times the number of possible categories for the column variable.

Filling in the cell counts

After you set up the table with the appropriate number of rows and columns, you need to fill in the appropriate numbers in each of the cells of the two-way table. The number in each cell of a two-way table is called the *cell count* for that cell. Of the four cells in the two-way table shown in Table 13-1, the upper-left cell represents the number of males who use their cellphones for personal calls. With the information you have in the

cellphone example, the cell count for this cell is 325. You also know that 427 females use their cellphones for personal calls, and this number goes into the lower-left cell.

To figure out the numbers in the remaining two cells, you do a bit of subtraction. You know from the information given that the total number of male cellphone users in the survey is 508. Each male either uses his cellphone for personal calls (falling into the *yes* group) or doesn't use it for personal calls (falling into the *no* group). Because 325 males fall into the *yes* group, and you have 508 males total, 183 males ($508 - 325 = 183$) don't use their cellphones for personal calls. This number is the cell count for the upper-right cell of the two-way table. Finally, because 508 females took the survey, and 427 of them use their cellphones for personal calls, you know that the rest of them ($508 - 427 = 81$) don't. Therefore, 81 is the cell count for the lower-right cell of the table. Table 13-2 shows the completed table for the cellphone user example, with the four cell counts filled in.

Table 13-2 Completed Two-Way Table for the Cellphone Data

	<i>Personal Calls: Yes</i>	<i>Personal Calls: No</i>
Males	325	$183 = (508 - 325)$
Females	427	$81 = (508 - 427)$



Just to save you a little time, if you have the total number in a group and the number of individuals who fall into one of the categories of the two-way table, you can determine the number falling into the remaining category by subtracting the total number in the group minus the number in the given category. You can complete this process for each remaining group in the table.

Making marginal totals

One of the most important characteristics of a two-way table is that it gives you easy access to all the pertinent totals. Because every two-way table is made up of rows and columns, you can imagine that the totals for each row and the totals for each column are important. Also, the grand total is important to know.

If you take a single row and add up all the cell counts in the cells of that row, you get a *marginal row total* for that row. Where does this marginal row total go on the table? You guessed it — out in the margin at the end of that row. You can find the marginal row totals for every row in the table and put them into the margins at the end of the rows. This group of marginal row totals for each row represents what statisticians call the *marginal distribution* for the row variable.



The marginal row totals should add up to the *grand total*, which is the total

number of individuals in the study. (The individuals may be people, cities, dogs, companies, and so on, depending on the scenario of the problem at hand.)

Similarly, if you take a single column and add up all the cell counts in the cells of that column, you get the *marginal column total* for that column. This number goes in the margin at the bottom of the column. Follow this pattern for each column in the table, and you have the marginal distribution for the column variable. Again, the sum of all the marginal column totals equals the grand total. The grand total is always located in the lower-right corner of the two-way table.

The marginal row total, marginal column totals, and the grand total for the cellphone example are shown in Table 13-3.

Table 13-3 Marginal and Grand Totals for the Cellphone Data

	<i>Personal Calls: Yes</i>	<i>Personal Calls: No</i>	<i>Marginal Row Totals</i>
<i>Males</i>	325	$183 = (508 - 325)$	508
<i>Females</i>	427	$81 = (508 - 427)$	508
<i>Marginal Column Totals</i>	752	264	1,016 (Grand Total)



The marginal row totals add the cell counts in each row; yet the marginal row totals show up as a column in the two-way table. This phenomenon occurs because when summing the cell counts in a row, you put the result in the margin at the end of the row, and when you do this for each row, you're stacking the row totals into a column. Similarly, the marginal column totals add the cell counts in each column; yet they show up as a row in the two-way table. Don't let this result be a source of confusion when you're trying to navigate or set up a two-way table. I recommend that you label your totals as marginal row, marginal column, or grand total to help keep it all clear.

Breaking Down the Probabilities

In the context of a two-way table, a percentage can be interpreted in one of two ways — in terms of a group or an individual. Regarding a group, a percentage represents the portion of the group that falls into a certain category. However, a percentage also represents the probability that an individual selected at random from the group falls into a certain category.

A two-way table gives you the opportunity to find many different kinds of probabilities, which help you to find the answers to different questions about your data or to look at the data another way. In this section, I cover the three most important types of probabilities found in a two-way table: marginal probabilities, joint probabilities, and conditional probabilities. (For more complete coverage of these types of probabilities,



When you find probabilities based on a sample, as you do in this chapter, you have to realize that those probabilities pertain to that sample only. They don't transfer automatically to the population being studied. For example, if you take a random sample of 1,000 adults and find that 55 percent of them watch reality TV, this study doesn't mean that 55 percent of all adults in the entire population watch reality TV. (The media makes this mistake every day.) You need to take into account the fact that sample results vary; in Chapters 14 and 15, you do just that. But this chapter zeros in on summarizing the information in your sample, which is the first step toward that end (but not the last step in terms of making conclusions about your corresponding population).

Marginal probabilities

A *marginal probability* makes a probability out of the marginal total, for either the rows or the columns. A marginal probability represents the proportion of the entire group that belongs in that single row or column category. Each marginal probability represents only one category for only one variable — it doesn't consider the other variable at all. In the cellphone example, you have four possible marginal probabilities (refer to Table 13-3):

- ✓ Marginal probability of female ($\frac{508}{1,016} = 0.50$), meaning that 50 percent of all the cellphone users in this sample were females
- ✓ Marginal probability of male ($\frac{508}{1,016} = 0.50$), meaning that 50 percent of all the cellphone users in this sample were males
- ✓ Marginal probability of using a cellphone for personal calls ($\frac{752}{1,016} = 0.74$), meaning that 74 percent of all cellphone users in this sample make personal calls with their cellphones
- ✓ Marginal probability of not using a cellphone for personal calls ($\frac{264}{1,016} = 0.26$), meaning that 26 percent of all the cellphone users in this sample don't make personal calls with their cellphones

Statisticians use shorthand notation for all probabilities. If you let M = male, F = female, Yes = personal cellphone use, and No = no personal cellphone use, then the preceding marginal probabilities are written as follows:

- ✓ $P(F) = 0.50$
- ✓ $P(M) = 0.50$
- ✓ $P(\text{Yes}) = 0.74$

✓ $P(\text{No}) = 0.26$

Notice that $P(F)$ and $P(M)$ add up to 1.00. This result is no coincidence because these two categories make up the entire gender variable. Similarly, $P(\text{Yes})$ and $P(\text{No})$ sum up to 1.00 because those choices are the only two for the personal cellphone use variable. Everyone has to be classified somewhere.



Be advised that some probabilities aren't useful in terms of discovering information about the population in general. For example, $P(F) = 0.50$ because the researchers determined ahead of time that they wanted exactly 508 females and exactly 508 males. The fact that 50 percent of the sample is female and 50 percent of the sample is male doesn't mean that in the entire population of cellphone users 50 percent are males and 50 percent are females. If you want to study what proportion of cellphone users are females and males, you need to take a combined sample instead of two separate ones and see how many males and females appear in the combined sample.

Joint probabilities

A *joint probability* gives the probability of the intersection of two categories, one from the row variable and one from the column variable. It's the probability that someone selected from the whole group has two particular characteristics at the same time. In other words, both characteristics happen jointly, or together. You find a joint probability by taking the cell count for those having both characteristics and dividing by the grand total.

Here are the four joint probabilities in the cellphone example:

- ✓ The probability that someone from the entire group is male and uses his cellphone for personal calls is $\frac{325}{1,016} = 0.32$, meaning that 32 percent of all the cellphone users in this sample are males using their cellphones for personal calls.
- ✓ The probability that someone from the entire group is male and doesn't use his cellphone for personal calls is $\frac{183}{1,016} = 0.18$.
- ✓ The probability that someone from the entire group is female and makes personal calls with her cellphone is $\frac{427}{1,016} = 0.42$.
- ✓ The probability that someone from the entire group is female and doesn't make personal calls with her cellphone is $\frac{81}{1,016} = 0.08$.

The notation for the joint probabilities listed is as follows, where + represents the intersection of the two categories listed:

- ✓ $P(M + \text{Yes}) = 0.32$
- ✓ $P(M + \text{No}) = 0.18$
- ✓ $P(F + \text{Yes}) = 0.42$
- ✓ $P(F + \text{No}) = 0.08$



The sum of all the joint probabilities for any two-way table should be 1.00, unless you have a little round-off error, which makes it very close to 1.00 but not exactly. The sum is 1.00 because everyone in the group is classified somewhere with respect to both variables. It's like dividing the entire group into four parts and showing which proportion falls into each part.

Conditional probabilities

A *conditional probability* is what you use to compare subgroups in the sample. In other words, if you want to break down the table further, you turn to a conditional probability. Each row has a conditional probability for each cell within the row, and each column has a conditional probability for each cell within that column.

Note: Because conditional probability is one of the sticking points for a lot of students, I spend extra time on it. My goal with this section is for you to have a good understanding of what a conditional probability really means and how you can use it in the real world (something many statistics textbooks neglect to mention, I have to say).

Figuring conditional probabilities

To find a conditional probability, you first look at a single row or column of the table that represents the known characteristic about the individuals. The marginal total for that row (column) now represents your new grand total, because this group becomes your entire universe when you examine it. Then take the cell counts from that row (column) and divide the sum by the marginal total for that row (column).

Consider the cellphone example in Table 13-3. Suppose you want to look at just the males who took the survey. The total number of males is 508. You can break down this group into two subgroups by using conditional probability: You can find the probability of using cellphones for personal calls (males only) and the probability of not using cellphones for personal calls (males only). Similarly, you can break down the females into those females who use cellphones for personal calls and those females who don't.

In the cellphone example, you have the following conditional probabilities when you break down the table by gender:

- ✓ The conditional probability that a male uses a cellphone for personal calls is

$$\frac{325}{508} = 0.64$$

- ✓ The conditional probability that a male doesn't use a cellphone for personal calls is $\frac{183}{508} = 0.36$.
- ✓ The conditional probability that a female uses a cellphone for personal calls is $\frac{427}{508} = 0.84$.
- ✓ The conditional probability that a female doesn't use a cellphone for personal calls is $\frac{81}{508} = 0.16$.

To interpret these results, you say that within this sample, if you're male, you're more likely than not to use your cellphone for personal calls (64 percent compared to 36 percent). However, the percentage of personal-call makers is higher for females (84 percent versus 16 percent).



Notice that for the males in the previous example, the two conditional probabilities (0.64 and 0.36) add up to 1.00. This is no coincidence. The males have been broken down by cellphone use for personal calls, and because everyone in the study is a cellphone user, each male has to be classified into one group or the other. Similarly, the two conditional probabilities for the females sum to 1.00.

Notation for conditional probabilities

You denote conditional probabilities with a straight vertical line that lists and separates the event that's known to have happened (what's given) and the event for which you want to find the probability. You can write the notation like this: $P(XX|XX)$. You place the given event to the right of the line and the event for which you want to find the probability to the left of the line. For example, suppose you know someone is female (F) and you want to find out the chance she's a Democrat (D). In this case, you're looking for $P(D|F)$. On the other hand, say you know a person is a Democrat and you want the probability that person is female — you're looking for $P(F|D)$.



The vertical line in the conditional probability notation isn't a division sign; it's just a line separating events A and B. Also, be careful of the order in which you place A and B into the conditional probability notation. In general, $P(A|B) \neq P(B|A)$.

Following is the notation used for the conditional probabilities in the cellphone example:

- ✓ **$P(\text{Yes}|\text{M}) = 0.64$** . You can say it this way: "The probability of Yes given Male is 0.64."

- ✓ **P(No | M) = 0.36.** In human terms, say, “The probability of No given Male is 0.36.”
- ✓ **P(Yes | F) = 0.84.** Say this one with gusto: “The probability of Yes given Female is 0.84.”
- ✓ **P(No | F) = 0.16.** You translate this notation by saying, “The probability of No given Female is 0.16.”



You can see that $P(\text{Yes} | \text{M}) + P(\text{No} | \text{M}) = 1.00$ because you’re breaking all males into two groups: those using cellphones for personal calls (Y) and those not (N). Notice, however, that $P(\text{Yes} | \text{M}) + P(\text{Yes} | \text{F})$ doesn’t sum to 1.00. In the first case, you’re looking only at the males, and in the second case, only at the females.

Comparing two groups with conditional probabilities

One of the most common questions regarding two categorical variables is this: Are they related? To answer this question, you compare their conditional probabilities.



To compare the conditional probabilities, follow these steps:

- 1. Take one variable and find the conditional probabilities based on the other variable.**
- 2. Repeat step one for each category of the first variable.**
- 3. Compare those conditional probabilities (you can even graph them for the two groups) and see whether they’re the same or different.**

If the conditional probabilities are the same for each group, the variables aren’t related in the sample. If they’re different, the variables are related in the sample.

- 4. Generalize the results to the entire population by using the sample results to draw a conclusion from the overall population involved by doing a Chi-square test (see Chapter 14).**

Revisiting the cellphone example from the previous section, you can ask specifically: Is personal use related to gender? You know that you want to compare cellphone use for males and females to find out whether use is related to gender. However, it’s very difficult to compare cell counts; for example, 325 males use their phones for personal calls, compared to 427 females. In fact, it’s impossible to compare these numbers without using some total for perspective. 325 out of what?



You have no way of comparing the cell counts in two groups without creating percentages (achieved by dividing each cell count by the appropriate total).

Percentages give you a means of comparing two numbers on equal terms. For example, suppose you gave a one-question opinion survey (yes, no, and no opinion) to a random sample of 1,099 people; 465 respondents said yes, 357 said no, and 277 had no opinion. To truly interpret this information, you're probably trying to compare these numbers to each other in your head. That's what percentages do for you. Showing the percentage in each group in a side-by-side fashion gives you a relative comparison of the groups with each other.

But first, you need to bring conditional probabilities into the mix. In the cellphone example, if you want the percentage of females who use their cellphones for personal calls, you take 427 divided by the total number of females (508) to get 84 percent. Similarly, to get the percentage of males who use their cellphones for personal calls, take the cell count (325) and divide it by that row total for males (508), which gives you 64 percent. This percentage is the conditional probability of using a cellphone for personal calls, given the person is male.

Now you're ready to compare the males and females by using conditional probabilities. Take the percentage of females who use their cellphones for personal calls and compare it to the percentage of males who use their cellphones for personal calls. By finding these conditional probabilities, you can easily compare the two groups and say that in this sample at least, more females use their cellphones (84 percent) for personal calls than men (64 percent).

Using graphs to display conditional probabilities

One way to highlight conditional probabilities as a tool for comparing two groups is to use graphs, such as a pie chart comparing the results of the other variable for each group or a bar chart comparing the results of the other variable for each group. (For more info on pie charts and bar charts, see my book *Statistics For Dummies* (Wiley) or your Stats I textbook.)

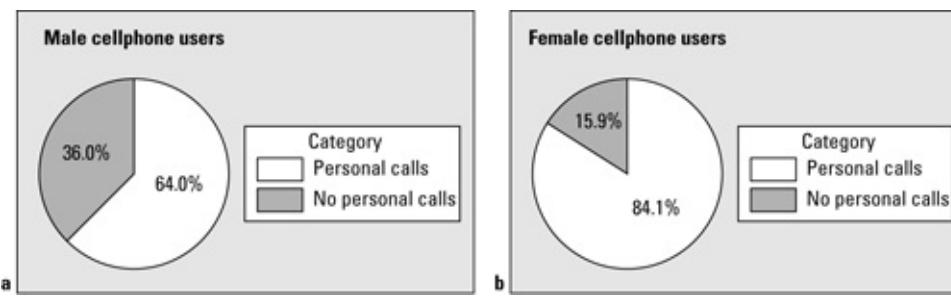


You may be wondering how close the two pie charts need to look (in terms of how close the slice amounts are for one pie compared to the other) in order to say the variables are independent. This question isn't one you can answer completely until you conduct a hypothesis test for the proportions themselves (see the Chi-square test in Chapter 14). For now, with respect to your sample data, if the difference in the appearance of the slices for the two graphs is enough that you would write a newspaper article about it, then go for dependence. Otherwise, conclude independence.

Figures 13-1a and 13-1b use two pie charts to compare cellphone use of males and females. Figure 13-1a shows the conditional distribution of cellphone use for (given) males. Figure 13-1b shows the conditional distribution of cellphone use for (given) females. A comparison of Figures 13-1a and 13-1b reveals that the slices for cellphone use

aren't equal (or even close) for males compared to females, meaning that gender and cellphone use for personal calls are dependent in this sample. This confirms the previous conclusions.

Figure 13-1:
Pie charts comparing male versus female personal cellphone use.



Another way you can make comparisons is to break down the two-way table by the column variable. (You don't always have to use the row variable for comparisons.) In the cellphone example (Table 13-3), you can compare the group of personal-call makers to the group of nonpersonal-call makers and see what percentage in each group is male and female. This type of comparison puts a different spin on the information because you're comparing the behaviors to each other in terms of gender.

With this new breakdown of the two-way table, you get the following:

- ✓ The conditional probability of being male, given you use your cellphone for personal calls, is $P(M | \text{Yes}) = \frac{325}{752} = 0.43$. **Note:** The denominator is 752, the total number of people who make personal calls with their cellphones.
- ✓ The conditional probability of being female, given you use your cellphone for personal calls, is $P(F | \text{Yes}) = \frac{427}{752} = 0.57$.

Again, these two probabilities add up to 1.00 because you're breaking down the personal-call makers according to gender (male or female). The conditional probabilities for the nonpersonal cellphone users are $P(M | \text{No}) = \frac{183}{264} = 0.69$ and $P(F | \text{No}) = \frac{81}{264} = 0.31$. These two probabilities also sum to 1.00 because you're breaking down the nonpersonal-call makers by gender (male and female).

The overall conclusions are similar to those found in the previous section, but the specific percentages and the interpretation are different. Interpreting the data this way, if you use your cellphone for personal calls, you're more likely to be female than male (57 percent compared to 43 percent). And if you don't use your cellphone to make personal calls, you're more likely to be male (69 percent compared to 31 percent).

What should you divide by? That is the question!

To get the correct answer for any probability in a two-way table, here's the trick: Always identify the group being examined. What's the probability "out of"? In the cellphone example (refer to Table 13-3),

- ✓ If you want the percentage of all users who are males using their phones for personal calls, you take the cell count 325 divided by 1,016, the grand total.

✓ If you want the percentage of males who are using their cellphones for personal calls, you take 325 divided by 508, the total number of males.

✓ If you want the percentage of personal-call makers who are male, you take 325 divided by 752, the total number of people who make personal calls with their cellphones.

In each of these three cases, the numerator is the same but the denominators are different, leading you to very different answers. Deciding which number to divide by is a very common source of confusion for people, and this trick can really give you an edge on keeping it straight.

Trying To Be Independent

Independence is a big deal in statistics. The term generally means that two items have outcomes whose probabilities don't affect each other. The items could be events A and B, variables x and y , survey results from two people selected at random from a population, and so on. If the outcomes of the two items do affect each other, statisticians call those two items *dependent* (or not independent). In this section, you check for and interpret independence of individual categories, one from each categorical variable in a sample, and you check for and interpret independence of two categorical variables in a sample.

Checking for independence between two categories

Statistics instructors often have students check to see whether two categories (one from a categorical variable x and the other from a categorical variable y) are independent. I prefer to just compare the two groups and talk about how similar or different the percentages are, broken down by another variable. However, to cover all the bases and make sure you can answer this very popular question, here's the official definition of independence, straight from the statistician's mouth: Two categories are *independent* if their joint probability equals the product of their marginal probabilities. The only caveat here is that neither of the categories can be completely empty.

For example, if being female is independent of being a Democrat, then $P(F + D) = P(F) * P(D)$, where D = Democrat and F = Female. So, to show that two categories are independent, find the joint probability and compare it to the product of the two marginal probabilities. If you get the same answer both times, the categories are independent. If not, then the categories are dependent.

You may be wondering: Don't all probabilities work this way, where the joint probability equals the product of the marginals? No, they don't. For example, if you draw a card from a standard 52-card deck, you get a red card with probability $\frac{1}{2}$. You draw a heart with probability $\frac{1}{4}$. The chance of drawing both a heart and a red card with one draw is still $\frac{1}{4}$.

(because all hearts are red).

However, the product of the individual probabilities for red and heart comes out to $\frac{1}{4} \cdot \frac{1}{2} = \frac{1}{8}$ which is not equal to $\frac{1}{4}$. This tells you that the categories “red” and “heart” aren’t independent (that is, they’re dependent).

Now the joint probability of a red two is $\frac{2}{52}$, or $\frac{1}{52}$. This equals the probability of a red card, $\frac{1}{2}$, times the probability of a two (because $\frac{1}{2} \cdot \frac{4}{52} = \frac{1}{26}$). This tells you that the categories “red” and “two” are independent.

Another way to check for independence is to compare the conditional probability to the marginal probability. Specifically, if you want to check whether being female is independent of being Democrat, check either of the following two situations (they’ll both work if the variables are independent):

- ✓ **Is $P(F|D) = P(F)$?** That is, if you know someone is a Democrat, does that affect the chance that they’ll also be female? If yes, F and D are independent. If not, F and D are dependent.
- ✓ **Is $P(D|F) = P(D)$?** This question is asking whether being female changes your chances of being a Democrat. If yes, D and F are independent. If not, D and F are dependent.



Is knowing that you’re in one category going to change the probability of being in another category? If so, the two categories aren’t independent. If knowing doesn’t affect the probability, then the two categories are independent.

Checking for independence between two variables



The previous section focuses on checking whether two specific categories are independent in a sample. If you want to extend this idea to showing that two entire categorical variables are independent, you must check the independence conditions for every combination of categories in those variables. All of them must work, or independence is lost. The first case where dependence is found between two categories means that the two variables are dependent. If you find that the first case shows independence, you must continue checking all the combinations before declaring independence.

Suppose a doctor’s office wants to know whether calling patients to confirm their appointments is related to whether they actually show up. The variables are x = called the patient (called or didn’t call) and y = patient showed up for his appointment (showed

or didn't show). Here are the four conditions that need to hold before you declare independence:

- ✓ $P(\text{showed}) = P(\text{showed} \mid \text{called})$
- ✓ $P(\text{showed}) = P(\text{showed} \mid \text{didn't call})$
- ✓ $P(\text{didn't show}) = P(\text{didn't show} \mid \text{called})$
- ✓ $P(\text{didn't show}) = P(\text{didn't show} \mid \text{didn't call})$

If any one of these conditions isn't met, you stop there and declare the two variables to be dependent in the sample. Only if all the conditions are met do you declare the two variables independent in the sample.

You can see the results of a sample of 100 randomly selected patients for this example scenario in Table 13-4.

Table 13-4 Confirmation Calls Related to Showing Up for the Appointment

	<i>Called</i>	<i>Didn't Call</i>	<i>Row Totals</i>
<i>Showed</i>	57	33	90
<i>Didn't Show</i>	3	7	10
<i>Column Totals</i>	60	40	100

Checking the conditions for independence, you can start at the first condition and check to see whether $P(\text{showed}) = P(\text{showed} \mid \text{called})$. From the last column of Table 13-4, you can see that $P(\text{showed})$ is equal to $\frac{90}{100} = 0.90$, or 90 percent. Next, look at the first column to find $P(\text{showed} \mid \text{called})$; this probability is $\frac{57}{60} = .95$ percent. Because these two probabilities aren't equal (although they're close), you say that showing up and calling first are dependent in this sample. You can also say that people come a little more often when you call them first. (To determine whether these sample results carry through to the population, which also takes care of the question of how close the probabilities need to be in order to conclude independence, see Chapter 14.)

Demystifying Simpson's Paradox

Simpson's Paradox is a phenomenon in which results appear to be in direct contradiction to one another, which can make even the best student's heart race. This situation can go unnoticed unless three variables (or more) are examined, in which case you organize the results into a *three-way table*, with columns within columns or rows within rows.

Simpson's Paradox is a favorite among statistics instructors (because it's so mystical and magical — and the numbers get so gooey and complex), but it's a nonfavorite among many students, mainly because of the following two reasons (in my opinion):

- ✓ Due to the way Simpson's Paradox is presented in most statistics courses, you can

easily get buried in the details and have no hope of seeing the big picture. Simpson's Paradox draws attention to a big problem in terms of interpreting data, and you need to understand the paradox fully in order to avoid it.

- Most textbooks do a good job of showing you examples of Simpson's Paradox, but they fall short in explaining why it occurs, so it just looks like smoke and mirrors. Some even neglect to explain the why part at all!

This section helps you get a handle on what Simpson's Paradox is, better understand why and how it happens, and know how to watch for it.

Experiencing Simpson's Paradox

Simpson's Paradox was discovered in 1951 by an American statistician named E. H. Simpson. He realized that if you analyze some data sets one way by breaking them down by two variables only, you can get one result, but when you break down the data further by a third variable, the results switch direction. That's why his result is called *Simpson's Paradox* — a paradox being an apparent contradiction in results.

Simpson's Paradox in action: Video games and the gender gap

The best way to sort through Simpson's Paradox is to watch it play out in an example and explain all the whys along the way. Suppose I'm interested in finding out who's better at playing video games, men or women. I watch males and females choose and play a variety of video games, and I record whether the player wins or loses. Suppose I record the results of 200 video games, as seen in Table 13-5. (Note that the females played 120 games, and the males played 80 games.)

Table 13-5 **Video Games Won and Lost for Males Versus Females**

	<i>Won</i>	<i>Lost</i>	<i>Marginal Row Totals</i>
<i>Males</i>	44	36	80
<i>Females</i>	84	36	120
<i>Marginal Column Totals</i>	128	72	200 (Grand Total)

Looking at Table 13-5, you see the proportion of males who won their video games, $P(\text{Won} \mid \text{Male})$, is $\frac{44}{80} = 0.55$. The proportion of females who won their video games, $P(\text{Won} \mid \text{Female})$, is $\frac{84}{120} = 0.70$. So overall, the females won more of their video games than the males did. Does this finding mean that women are better than men at video games in general in the sample?

Not so fast, my friend. Notice that the people in the study were allowed to choose the video games they played. This factor blows the study wide open. Suppose females and males choose different types of video games: Can this affect the results? The answer may

be yes. Considering other variables that could be related to the results but weren't included in the original study (or at least not in the original data analysis) is important. These additional variables that cloud the results are called *lurking variables*.

Factoring in difficulty level

Many people may expect the video game results from the previous section to be turned around to indicate that men are better at playing video games than women. According to the research, men spend more time playing video games, on average, and are by far the primary purchasers of video games, compared to women. So what explains the eyebrow-raising results in this study? Is there another possible explanation? Is important information missing that's relevant to this case?

One of the variables that wasn't considered when I made Table 13-5 was the difficulty level of the video game being played. Suppose I go back and include the difficulty level of the chosen game each time, along with each result (won or lost). Level one indicates easy video games, comparable to the level of Ms. Pac Man (games that are my speed), and level two means more challenging video games (like war games or sophisticated strategy games).

Table 13-6 represents the results with the addition of this new information on difficulty level of games played. You have three variables now: level of difficulty (one or two), gender (male or female), and outcome (won or lost). That makes Table 13-6 a three-way table.

Table 13-6 A Three-Way Table for Gender, Game Level, and Game Outcome

Level-One Games		Level-Two Games	
	Won	Lost	Won
Males	9	1	35
Females	72	18	12

Note in Table 13-6 that the number of level-one video games chosen was $9 + 1 + 72 + 18 = 100$, and the number of level-two video games chosen was $35 + 35 + 12 + 18 = 100$. In order to reevaluate the data based on the game level information, you need to look at who chose which level of game. The next section probes this very issue.

Comparing success rates with conditional probabilities

To compare the success rates for males versus females using Table 13-6, you can figure out the appropriate conditional probabilities, first for level-one games and then for level-two games.

For level-one games (only), the conditional probability of winning given male is $P(\text{Won}|\text{Male}) = \frac{9}{10} = 0.90$. So for the level-one games, males won 90 percent of the games they played. For level-one games, the percentage of games won by the females is

$P(\text{Won}|\text{Female}) = \frac{72}{90} = 0.80$, or 80 percent. These results mean that at level one, the males did 10 percent better than the females at winning their games. But this percentage appears to contradict the results found in Table 13-5. (Just wait — the contradictions don't end here!)

Now figure the conditional probabilities for the level-two video games won.

For the men, the percentage of males winning level-two games was $\frac{35}{70} = 0.50$, or 50 percent. For the ladies, the percentage of women winning level-two games was $\frac{12}{30} = 0.40$, or 40 percent. Once again, the males outdid the females!

Step back and think about this scenario for a minute. Table 13-5 shows that females won a higher percentage of the video games they played overall. But Table 13-6 shows that males won more of the level-one games and more of the level-two games. What's going on? No need to check your math. No mistakes were made — no tricks were pulled. This inconsistency in results happens in real life from time to time in situations where an important third variable is left out of a study, a situation aptly named *Simpson's Paradox*. (See why it's called a paradox?)

Figuring out why Simpson's Paradox occurs



Lurking variables are the underlying cause of Simpson's Paradox. A *lurking variable* is a third variable that's related to each of the other two variables and can affect the results if not accounted for.

In the video game example, when you look at the video game outcomes (won or lost) broken down by gender only (Table 13-5), females won a higher percentage of their overall games than males (70 percent overall winning percentage for females compared to 55 percent overall winning for males). Yet, when you split up the results by the level of the video game (level one or level two; see Table 13-6), the results reverse themselves, and you see that males did better than females on the level-one games (90 percent compared to 80 percent), and males also did better on the level-two games (50 percent compared to 40 percent).

To see why this seemingly impossible result happens, take a look at the marginal row *probabilities* versus the marginal row *totals* for the level-one games in Table 13-6. The percentage of times a male won when he played an easy video game was 90 percent. However, males chose level-one video games only 10 times out of 80 total level-one games played by men. That's only 12.5 percent.

To break this idea down further, the males' nonstellar performance on the challenging video games (50 percent — but still better than the females) coupled with the fact that the males chose challenging video games 87.5 percent of the time (that's 70 out of 80

times) really brought down their overall winning percentage (55 percent). And even though the men did really well on the level-one video games, they didn't play many of them (compared to the females), so their high winning percentage on level-one video games (90 percent) didn't count much toward their overall winning percentage.

Meanwhile, in Table 13-6, you see that females chose level-one video games 90 times (out of 120). Even though the females only won 72 out of the 90 games (80 percent, a lower percentage than the males, who won 9 out of 10 of their games), they chose to play many more level-one games, therefore boosting their overall winning percentage.

Now the opposite situation happens when you look at the level-two video games in Table 13-6. The males chose the harder video games 70 times (out of 80), while the females only chose the harder ones only 30 times out of 120. The males did better than the females on level-two video games (winning 50 percent of them versus 40 percent for the females). However, level-two video games are harder to win than level-one video games. This factor means that the males' winning percentage on level-two video games, being only 50 percent, doesn't contribute much to their overall winning percentage. However, the low winning percentage for females on level-two video games doesn't hurt them much, because they didn't play many level-two video games.



The bottom line is that the occurrence or nonoccurrence of Simpson's Paradox is a matter of weights. In the overall totals from Table 13-5, the males don't look as good as the females. But when you add in the difficulty of the games, you see that most of the males' wins came from harder games (which have a lower winning percentage). The females played many more of the easier games on average, and easy games carry a higher chance of winning no matter who plays them. So it all boils down to this: Which games did the males choose to play, and which games did the females choose to play? The males chose harder games, which contributed in a negative way to their overall winning percentage and made the females look better than they actually were.

Keeping one eye open for Simpson's Paradox

Simpson's Paradox shows you the importance of including data about possible lurking variables when attempting to look at relationships between categorical variables.

Level of game wasn't included in the original summary, Table 13-5, but it should have been included because it's a variable that affected the results. Level of game, in this case, was the lurking variable. More men chose to play the more difficult games, which are harder to win, thereby lowering their overall success rate.



You can avoid Simpson's Paradox by making sure that obvious lurking variables

are included in a study; that way, when you look at the data you get the relationships right the first time and there's a lower chance of reversing the results. And, as with all other statistical results, if it looks too good to be true or too simple to be correct, it probably is! Beware of someone who tried to oversimplify any result. While three-way tables are a little more difficult to examine, they're often worth using.

Being Independent Enough for the Chi-Square Test

In This Chapter

- ▶ Testing for independence in the population (not just the sample)
 - ▶ Using the Chi-square distribution
 - ▶ Discovering the connection between the Z-test and the Chi-square test
-

You've seen these hasty judgments before — people who collect one sample of data and try to use it to make conclusions about the whole population. When it comes to two categorical variables (where data fall into categories and don't represent measurements), the problem seems to be even more widespread.

For example, a TV news show finds that out of 1,000 presidential voters, 200 females are voting Republican, 300 females are voting Democrat, 300 males are voting Republican, and 200 males are voting Democrat. The news anchor shows the data and then states that 30 percent ($300 \div 1,000$) of all presidential voters are females voting Democrat (and so on for the other counts).

This conclusion is misleading. It's true that in this sample of 1,000 voters, 30 percent of them are females voting Democrat. However, this result doesn't automatically mean that 30 percent of the entire population of voters is females voting Democrat. Results change from sample to sample.

In this chapter, you see how to move beyond just summarizing the sample results from a two-way table (discussed in Chapter 13) to using those results in a hypothesis test to make conclusions about an entire population. This process requires a new probability distribution called the *Chi-square distribution*. You also find out how to answer a very popular question among researchers: Are these two categorical variables independent (not related to each other) in the entire population?

The Chi-square Test for Independence

Looking for relationships between variables is one of the most common reasons for collecting data. Looking at one variable at a time usually doesn't cut it. The methods used to analyze data for relationships are different depending on the type of data collected. If the two variables are quantitative (for example, study time and exam score), you use correlation and regression (see Chapter 4). If the two variables are categorical (for example, gender and political affiliation), you use a Chi-square test to examine relationships. In this section, you see how to use a Chi-square test to look for relationships between two categorical variables.



If two categorical variables don't have a relationship, they're deemed to be *independent*. If they do have a relationship, they're called *dependent variables*. Many folks get confused by these terms, so it's important to be clear about the distinction right up front.

To test whether two categorical variables are independent, you need a Chi-square test. The steps for the Chi-square test follow. (Minitab can conduct this test for you, from step three on down.)

1. Collect your data, and summarize it in a two-way table.

These numbers represent the observed cell counts. (For more on two-way tables, see Chapter 13.)

2. Set up your null hypothesis, H_0 : Variables are independent; and the alternative hypothesis, H_a : Variables are dependent.

3. Calculate the expected cell counts under the assumption of independence.

The expected cell count for a cell is the row total times the column total divided by the grand total.

4. Check the conditions of the Chi-square test before proceeding; each expected cell count must be greater than or equal to five.

5. Figure the Chi-square test statistic.

This statistic finds the observed cell count minus the expected cell count, squares the difference, and divides it by the expected cell count. Do these steps for each cell, and then add them all up.

6. Look up your test statistic on the Chi-square table (Table A-3 in the appendix) and find the p -value (or one that's close).

7. If your result is less than your predetermined cutoff (the α level), usually 0.05, reject H_0 and conclude dependence of the two variables.

If your result is greater than the α level, fail to reject H_0 ; the variables can't be deemed dependent.



To conduct a Chi-square test in Minitab, enter your data in the spreadsheet exactly as it appears in your two-way table (see Chapter 13 for setting up a two-way table for categorical data). Go to Stat>Tables>Chi-Square Test. Click on the two variable names in the left-hand box corresponding to your column variables in the spreadsheet. They appear in the box labeled Columns Contained in the Table. Then click OK.

Collecting and organizing the data

The first step in any data analysis is collecting your data. In the case of two categorical

variables, you collect data on the two variables at the same time for each individual in the study.

A survey conducted by American Demographics asked men and women about the color of their next house. The results showed that 36 percent of the men wanted to paint their houses white, and 25 percent of the women wanted to paint their houses white. Keeping the data together in pairs (for example: male, white paint; female, nonwhite paint), you organize them into a two-way table where the rows represent the categories of one categorical variable (males and females for gender) and the columns represent the categories of the other categorical variable (white paint and nonwhite paint). Table 14-1 contains the results from a sample of 1,000 people (500 men and 500 women).

Table 14-1 **Gender and House Paint Color Preference:**
Observed Cell Counts

	<i>White Paint</i>	<i>Nonwhite Paint</i>	<i>Marginal Row Totals</i>
<i>Men</i>	180	320	500
<i>Women</i>	125	375	500
<i>Marginal Column Totals</i>	305	695	1,000 (Grand Total)

The *marginal row totals* represent the total number in each row; the *marginal column totals* represent the total number in each column. (See Chapter 13 for more information on row and column marginal totals.)

Notice that of the males, the percentage that wants to paint the house white is $180 \div 500 = 0.36$, or 36 percent, as stated previously. And the percentage of females that wants to paint the house white is $125 \div 500 = 0.25$, or 25 percent. (Both of these percentages represent conditional probabilities as explained in Chapter 13.)

The American Demographics report concluded from this data that “. . . men and women generally agree on exterior house paint colors; the main exception being the top male choice, white (36 percent would paint their next house white versus 25 percent of women).” This type of conclusion is commonly formed, but it’s an overgeneralization of the results at this point.

You know that in this sample, more men wanted to paint their houses white than women, but is 180 really that different from 125 when you’re dealing with a sample size of 1,000 people whose results will vary the next time you do the survey? How do you know these results carry over to the population of all men and women? That question can’t be answered without a formal statistical procedure called a *hypothesis test* (see Chapter 3 for the basics of hypothesis tests).

To show that men and women in the population differ according to favorite house color, first note that you have two categorical variables:

- ✓ Gender (male or female)

- ✓ Paint color (white or nonwhite)



Making conclusions about the population based on the sample (observed) data in a two-way table is taking too big of a leap. You need to conduct a Chi-square test in order to broaden your conclusions to the entire population. The media, and even some researchers, can get into trouble by ignoring the fact that sample results vary. Stopping with the sample results only and going merrily on your way can lead to conclusions that others can't confirm when they take new samples.



You keep the connection between the two pieces of information by organizing the data into one two-way table versus two individual tables — one for gender and one for house-paint preference. With one two-way table, you can look at the relationship between the two variables. (For the full details on organizing and interpreting the results from a two-way table, see Chapter 13.)

Determining the hypotheses

Every hypothesis test (whether it be a Chi-square test or some other test) has two hypotheses:

- ✓ **Null hypothesis:** You have to believe this unless someone shows you otherwise. The notation for this hypothesis is H_0 .
- ✓ **Alternative hypothesis:** You want to conclude this in the event that you can't support the null hypothesis anymore. The notation for this hypothesis is H_a .

In the case where you're testing for the independence of two categorical variables, the null hypothesis is when no relationship exists between them. In other words, they're independent. The alternative hypothesis is when the two variables are related, or dependent.

For the paint color preference example from the previous section, you write H_0 : Gender and paint color preference are independent versus H_a : Gender and paint color preference are dependent. And there you have it — step two of the Chi-square test.

For a quick review of hypothesis testing, turn to Chapter 3. For a full discussion of the topic, see my other book *Statistics For Dummies* (Wiley) or your Stats I textbook.

Figuring expected cell counts

When you've collected your data and set up your two-way table (for example, see Table 14-1), you already know what the observed values are for each cell in the table. Now you

need something to compare them to. You’re ready for step three of the Chi-square test — finding expected cell counts.

The null hypothesis says that the two variables x and y are independent. That’s the same as saying x and y have no relationship. Assuming independence, you can determine which numbers should be in each cell of the table by using a formula for what’s called the *expected cell counts*. (Each individual square in a two-way table is called a *cell*, and the number that falls into each cell is called the *cell count*; see Chapter 13.)

Table 14-1 shows the observed cell counts from the gender and paint color preference example. To find the expected cell counts you take the row total times the column total divided by the grand total, and do this for each cell in the table. Table 14-2 shows the calculations for the expected cell counts for the gender and paint color preference data.

Table 14-2 **Gender and House Paint Color Preference:**
Expected Cell Counts

	<i>White Paint</i>	<i>Nonwhite Paint</i>	<i>Marginal Row Totals</i>
<i>Men</i>	$(500 * 305) \div 1000$ = 152.5	$(500 * 695) \div 1000$ = 347.5	500
<i>Women</i>	$(500 * 305) \div 1000$ = 152.5	$(500 * 695) \div 1000$ = 347.5	500
<i>Marginal Column Totals</i>	305	695	1000 (Grand Total)

Next you compare the observed cell counts in Table 14-1 to the expected cell counts in Table 14-2 by looking at their differences. The differences between the observed and expected cell counts shown in these tables are the following:

$$180 - 152.5 = 27.5$$

$$320 - 347.5 = -27.5$$

$$125 - 152.5 = -27.5$$

$$375 - 347.5 = 27.5$$

Next you do a Chi-square test for independence (see Chapter 15) to determine whether the differences found in the sample between the observed and expected cell counts are simply due to chance, or whether they carry through to the population.



Under independence, you conclude there is not a significant difference between what you observed and what you expected.

Checking the conditions for the test



Step four of the Chi-square test is checking conditions. The Chi-square test has one main condition that must be met in order to test for independence on a two-way table: The expected count for each cell must be at least five — that is, greater than or equal to five. Expected cell counts that fall below five aren't reliable in terms of the variability that can take place.

In the gender and paint color preference example, Table 14-2 shows that all the expected cell counts are at least five, so the conditions of the Chi-square test are met.



If you're analyzing data and you find that your data set doesn't meet the expected cell count of at least five for one or more cells, you can combine some of your rows and/or columns. This combination makes your table smaller, but it increases the cell counts for the cells that you do have, which helps you meet the condition.

Calculating the Chi-square test statistic

Every hypothesis test uses data to make the decision about whether or not to reject H_0 in favor of H_a . In the case of testing for independence in a two-way table, you use a hypothesis test based on the Chi-square test statistic. In the following sections, you can see the steps for calculating and interpreting the Chi-square test statistic, which is step five of the Chi-square test.

Working out the formula

A major component of the Chi-square test statistic is the expected cell count for each cell in the table. The formula for finding the expected cell count, e_{ij} , for the cell in row i ,

$$\text{column } j \text{ is } e_{ij} = \frac{\text{row } i \text{ total} * \text{column } j \text{ total}}{\text{grand total}}.$$

Note that the values of i and j vary for each cell in the table. In a two-way table, the upper-left cell of the table is in row one, column one. The cell in the upper-right corner is in row one, column two. The cell in the lower-left corner is in row two, column one, and the lower-right cell is in row two, column two.

The formula for the Chi-square test statistic is $\chi^2 = \sum_i \sum_j \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$, where o_{ij} is the observed cell count for the cell in row i , column j , and e_{ij} is the expected cell count for the cell in row i , column j .



When you calculate the expected cell count for some cells, you typically get a

number that has some digits after the decimal point (in other words, the number isn't a whole number). Don't round this number off, despite the temptation to do so. This expected cell count is actually an overall-average expected value, so keep the count as it is, with decimal included.

Calculating the test statistic

Here are the major steps of how to calculate the Chi-square test statistic for independence (Minitab does these steps for you as well):

- 1. Subtract the observed cell count from the expected cell count for the upper-left cell in the table.**
- 2. Square the result from step one to make the number positive.**
- 3. Divide the result from step two by the expected cell count.**
- 4. Repeat this process for all the cells in the table, and add up all the results to get the Chi-square test statistic.**



The reason you divide by the expected cell count in the Chi-square test statistic is to account for cell-count sizes. If you expect a big cell count, say 100, and are off by only 5 for the observed count of that cell, that difference shouldn't count as much as if you expected a small cell count (like 10) and the observed cell count was off by 5. Dividing by the expected cell count puts a more fair weight on the differences that go into the Chi-square test statistic.



To perform a Chi-square test in Minitab, you have to first enter the raw data (the data on each person) in two columns. The first column contains the values of the first variable in your data set. (For example, if your first variable is gender, go down the first Minitab column entering the gender of each person.) Then enter the data from your second variable in the second column, where each row represents a single person in the data set. (If your second variable is house paint color preference, for example, enter each person's paint color preference in column two, keeping the data from each person together in each row.) Go to Stat>Tables>Cross-tabulation and Chi-square.

Now Minitab needs to know which is your row variable and which is your column variable in your table. On the left-hand side, click on the variable that you want to represent the rows of your two-way table (you may click on the first variable). Click Select, and the variable name appears in the row variable portion of the table on the right. Now find the column variable blank on the right-hand side and click on it. Go to the left-hand side and click on the name of your second variable. Click Select. Then click on the Chi-square button and choose Chi-square analysis by checking the box. If you want the expected cell counts included, check that box also. Then click OK. Finally, click OK

again to clear all the windows.

Picking through the output

The Minitab output for the Chi-square analysis for the gender and house paint color preference example (from Table 14-1) is shown in Figure 14-1. You can pick out quite a few numbers from the output in Figure 14-1 that are especially important. The following three numbers are listed in each cell:

- ✓ The first (top) number is the observed cell count for that cell; this matches the observed cell count for each cell shown in Table 14-1. (Notice that the marginal row and column totals of Figure 14-1 also match those from Table 14-1.)
- ✓ The second number in each cell of Figure 14-1 is the expected cell count for that cell; you find it by taking the row total times the column total divided by the grand total (see the section “Figuring expected cell counts”). For example, the expected cell count for the upper-left cell (males who prefer white house paint) is $(500 * 305) / 1,000 = 152.50$.
- ✓ The third number in each cell of Figure 14-1 is that part of the Chi-square test statistic that comes from that cell. (See steps one through three of the previous section “Working out the formula.”) The sum of the third numbers in each cell equals the value of the Chi-square statistic listed in the last line of the output. (For the house paint color preference example, the Chi-square test statistic is 14.27.)

Figure 14-1:
Minitab
output for the
house paint
color
preference
data.

Chi-Square Test: Gender, House-Paint Preference				
Expected counts are printed below observed counts				
Chi-Square contributions are printed below expected counts				
	White Paint	Nonwhite Paint	Total	
M	180	320	500	
	152.50	347.50		
	4.959	2.176		
F	125	375	500	
	152.50	347.50		
	4.959	2.176		
Total	305	695	1000	
Chi-Sq = 14.271, DF = 1, P-Value = 0.000				

Finding your results on the Chi-square table

The only way to make an assessment about your Chi-square test statistic is to compare it to all the possible Chi-square test statistics you would get if you had a two-way table with the same row and column totals, yet you distributed the numbers in the cells in every way possible. (You can do that in your sleep, right?) Some resulting tables give large Chi-square test statistics, and some give small Chi-square test statistics.

Putting all these Chi-square test statistics together gives you what's called a *Chi-square distribution*. You find your particular test statistic on that distribution (step six of the Chi-square test), and see where it stands compared to the rest.

If your test statistic is large enough that it appears way out on the right tail of the Chi-square distribution (boldly going where no test statistic has gone before), you reject H_0 and conclude the two variables are not independent. If the test statistic isn't that far out, you can't reject H_0 .

In the next sections, you find out more about the Chi-square distribution and how it behaves, so you can make a decision about the independence of your two variables based on your Chi-square statistic.

Determining degrees of freedom

Each type of two-way table has its own Chi-square distribution, depending on the number of rows and columns it has, and each Chi-square distribution is identified by its *degrees of freedom*.

In general, a two-way table with r rows and c columns uses a Chi-square distribution with $(r - 1) * (c - 1)$ degrees of freedom. A two-way table with two rows and two columns uses a Chi-square distribution with one degree of freedom. Notice that $1 = (2 - 1) * (2 - 1)$. A two-way table with three rows and two columns uses a Chi-square distribution with $(3 - 1) * (2 - 1) = 2$ degrees of freedom.



Understanding *why* degrees of freedom are calculated this way is likely to be beyond the scope of your statistics class. But if you really want to know, the degrees of freedom represents the number of cells in the table that are flexible, or free, given all the marginal row and column totals.

For example, suppose that a two-way table has all row and column totals equal to 100 and the upper-left cell is 70. Then the upper-right cell must be 100 (row total) – $30 = 70$. Because the column one total is 100, and the upper-left cell count is 70, the lower-left cell count must be $100 - 70 = 30$. Similarly, the lower-right cell count must be 70.

So you have only one free cell in a two-way table after you have the marginal totals set up. That's why the degree of freedom for a two-way table is 1. In general, you always lose one row and one column because of knowing the marginal totals. That's because the last row and column values can be calculated through subtraction. That's where the formula $(r - 1) * (c - 1)$ comes from. (That's more than you wanted to know, isn't it?)

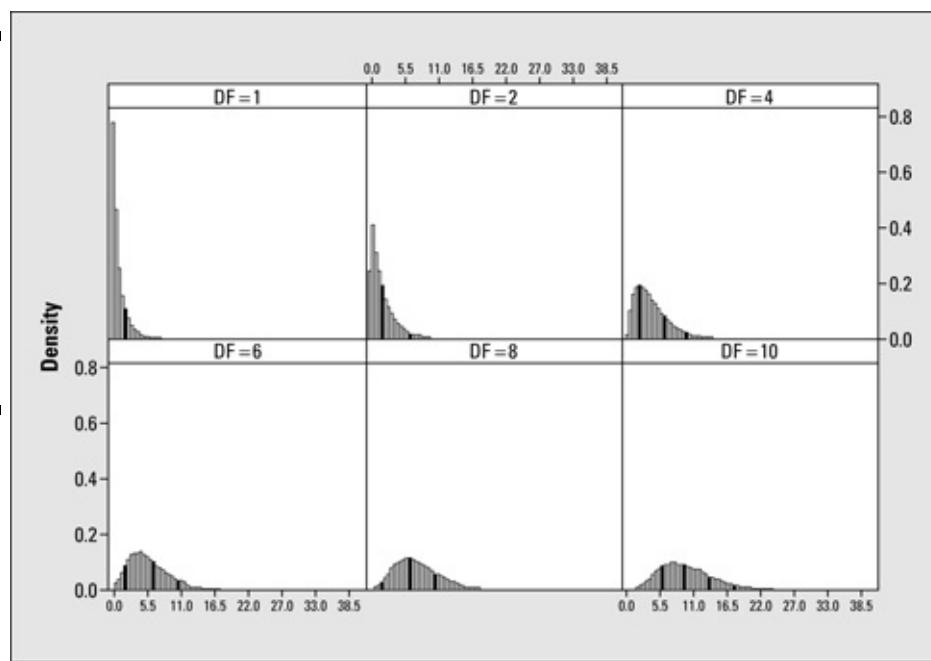
Discovering how Chi-square distributions behave

Figure 14-2 shows pictures of Chi-square distributions with 1, 2, 4, 6, 8, and 10 degrees of

freedom, respectively. Here are some important points to keep in mind about Chi-square distributions:

- ✓ For 1 degree of freedom, the distribution looks like a hyperbola (see Figure 14-2, top left); for more than 1 degree of freedom, it looks like a mound that has a long right tail (see Figure 14-2, lower right).
- ✓ All the values are greater than or equal to zero.
- ✓ The shape is always skewed to the right (tail going off to the right).
- ✓ As the number of degrees of freedom increases, the mean (the overall average) increases (moves to the right) and the variances increase (resulting in more spread).
- ✓ No matter what the degree of freedom is, the values on the Chi-square distribution (known as the *density*) approach zero for increasingly larger Chi-square values. That means that larger and larger Chi-square values are less and less likely to happen.

Figure 14-2:
Chi-square distributions with 1, 2, 4, 6, 8, and 10 degrees of freedom (moving from upper left to lower right).



Using the Chi-square table

After you find your Chi-square test statistic and its degrees of freedom, you want to determine how large your statistic is, relative to its corresponding distribution. (You're now venturing into step seven of the Chi-square test.)

If you think about it graphically, you want to find the probability of being beyond (getting a larger number than) your test statistic. If that probability is small, your Chi-square test statistic is something unusual — it's out there — and you can reject H_0 . You then conclude that your two variables are not independent (they're related somehow).



In case you’re following along at home, the Chi-square test statistic for the independent data from Table 14-2 is zero because the observed cell counts are equal to the expected cell counts for each cell, and their differences are always equal to zero. (This result never happens in real life!) This scenario represents a *perfectly independent* situation and results in the smallest possible value of a Chi-square test statistic.

If the probability of being to the right of your Chi-square test statistic (on a graph) isn’t small enough, you don’t have enough evidence to reject H_0 . You then stick with H_0 ; you can’t reject it. You conclude that your two variables are independent (unrelated).



How small of a probability do you need to reject H_0 ? For most hypothesis tests, statisticians generally use 0.05 as the cutoff. (For more information on cutoff values, also known as α levels, flip to Chapter 3, or check out my other book *Statistics For Dummies*.)

Your job now is to find the probability of being beyond your Chi-square test statistic on the corresponding Chi-square distribution with $(r - 1) * (c - 1)$ degrees of freedom. Each Chi-square distribution is different, and because the number of possible degrees of freedom is infinite, showing every single value of every Chi-square distribution isn’t possible.

In the Chi-square table in (Table A-3 in the appendix), you see some of the most important values on each Chi-square distribution with degrees of freedom from 1 to 50.

To use the Chi-square table, you find the row that represents your degrees of freedom (abbreviated df). Move across that row until you reach the value closest to your Chi-square test statistic, without going over. (It’s like a game show where you’re trying to win the showcase by guessing the price.)

Then go to the top of the column you’re in. That number represents the area to the right (above) of the Chi-square test statistic you saw in the table. The area above your particular Chi-square test statistic is less than or equal to this number. This result is the approximate p -value of your Chi-square test.

In the house paint color preference example (see Figure 14-1), the Chi-square test statistic is 14.27. You have $(2 - 1) * (2 - 1) = 1$ degree of freedom. In the Chi-square table, go to the row for $df = 1$, and go across to the number closest to 14.27 (without going over), which is 7.88.

Drawing your conclusions

You have two alternative ways to draw conclusions from the Chi-square test statistic. You can look up your test statistic on the Chi-square table and see the probability of being greater than that. This method is known as *approximating the p-value*. (The *p*-value of a test statistic is the probability of being at or beyond your test statistic on the distribution to which the test statistic is being compared — in this case, the Chi-square distribution.) Or you can have the computer calculate the exact *p*-value for your test. (For a quick review of *p*-values and α levels, turn to Chapter 3. For a full review of these topics, see my other book *Statistics For Dummies*.)

Before you do anything though, set your α , the cutoff probability for your *p*-value, in advance. If your *p*-value is less than your α level, reject H_0 . If it's more, you can't reject H_0 .

Approximating p-value from the table

For the house paint color preference example (see Figure 14-1), the Chi-square test statistic is 14.27 with $(2 - 1) * (2 - 1) = 1$ df (degree of freedom). The closest number in row one of the Chi-square table (see Table A-3 in the appendix), without going over, is 7.88 (in the last column).

The number at the top of that column is 0.005. This number is less than your typical α level of 0.05, so you reject H_0 . You know that your *p*-value is less than 0.005 because your test statistic was more than 7.88. In other words, if 7.88 is the minimum evidence you need to reject H_0 , you have more evidence than that with a value of 14.28. More evidence against H_0 means a smaller *p*-value.

However, because Chi-square tables in general only give a few values for each Chi-square distribution, the best you can say using this table is that your *p*-value for this test is less than 0.005.

Here's the big news: Because your *p*-value is less than 0.05, you can conclude based on this data that gender and house paint color preference are likely to be related in the population (dependent), like the American Demographics Survey said (quoted at the beginning of this chapter). Only now, you have a formal statistical analysis that says this result found in the sample is also likely to occur in the entire population. This statement is much stronger!



If your data shows you can reject H_0 , you only know at that point that the two variables have some relationship. The Chi-square test statistic doesn't tell you what that relationship is. In order to explore the relationship between the two variables, you find the conditional probabilities in your two-way table (see Chapter 13). You can use those results to give you some ideas as to what may be happening in the population.

For the gender and house paint color preference example, because paint color preference is related to gender, you can examine the relationship further by comparing the male versus female paint color preferences and describing how they're different. Start by finding the percentage of men that prefer white houses, which comes out to $180 \div 500 = 0.36$, or 36 percent, calculated from Table 14-1. Now compare this result to the percentage of women who prefer white houses: $125 \div 500 = 0.25$, or 25 percent. You can now conclude that in this population (not just the sample) men prefer white houses more than women do. Hence, gender and house paint color preference are dependent.



Dependent variables affect each other's outcomes, or cell counts. If the cell counts you actually observe from the sample data won't match the expected cell counts under H_0 : The variables are independent, you conclude that the dependence relationship you found in the sample data carries over to the population. In other words, big differences between observed and expected cell counts mean that the variables are dependent.

Extracting the *p*-value from computer output

After Minitab calculates the test statistic for you, it reports the exact *p*-value for your hypothesis test. The *p*-value measures the likelihood that your results were found just by chance while H_0 is still true. It tells you how much strength you have against H_0 . If the *p*-value is 0.001, for example, you have much more strength against H_0 than if the *p*-value, say, is 0.10.

Looking at the Minitab output for the gender-paint color preference data in Figure 14-1, the *p*-value is reported to be 0.000. This means that the *p*-value is smaller than 0.001; for example, it may be 0.0009. That's a very small *p*-value! (Minitab only reports results to three decimal points, which is typical of many statistical software packages.)



I've seen situations where people get a result that isn't quite what they want (like a *p*-value of 0.068), and so they do some tweaking to get what they want. They change their α level from 0.05 to 0.10 after the fact. This change makes the *p*-value less than the α level, and they feel they can reject H_0 and say that a relationship exists.

But what's wrong with this picture? They changed the α after they looked at the data, which isn't allowed. That's like changing your bet in blackjack after you find out what the dealer's cards look like. (Tempting, but a serious no-no.) Always be wary of large α levels, and make sure that you always choose your α before collecting any data — and stick to it.

The good news is that when *p*-values are reported, anyone reading them can make his or her own conclusion; no cut-and-dry rejection and acceptance region is set in stone. But

setting an α level once and then changing it after the fact to get a better conclusion is never good!

Putting the Chi-square to the test

If two variables turn out to be dependent, you can describe the relationship between them. But if two variables are independent, the results are the same for each group being compared. The following example illustrates this idea.

There has been much speculation and debate as to whether cellphone use should be banned while driving. You're interested in Americans' opinions on this issue, but you also suspect that the results may differ by gender. You decide to do a Chi-square test for independence to see if your theory plays out. Table 14-3 shows a two-way table of observed data from 60 men and 60 women regarding whether they agree with the policy (banning cellphone use while driving) or not. From Table 14-3 you see that $12 \div 60 = 20$ percent of men agree with the policy of banning cellphones while driving, compared to $9 \div 60 = 15$ percent of women. You see these percentages are different, but is this enough to say that gender and opinion on this issue are dependent? Only a Chi-square test for independence can help you decide.

Table 14-3 Gender and Opinion on Cellphone Ban:
Observed Cell Counts

	Agree with Cellphone Ban	Disagree with Cellphone Ban	Marginal Row Totals
Men	12	48	60
Women	9	51	60
Marginal Column Totals	21	99	120 (Grand Total)

Table 14-4 shows the expected cell counts under H_0 , along with their calculations.

Table 14-4 Gender and Opinion on Cellphone Ban:
Expected Cell Counts

	Agree with Cellphone Ban	Disagree with Cellphone Ban	Marginal Row Totals
Men	$(60 * 21) \div 120$ $= 10.5$	$(60 * 99) \div 120$ $= 49.5$	60
Women	$(60 * 21) \div 120$ $= 10.5$	$(60 * 99) \div 120$ $= 49.5$	60
Marginal Column Totals	21	99	120 (Grand Total)

Running a Chi-square test in Minitab for this data, the degrees of freedom equals $(2 - 1) * (2 - 1) = 1$; the Chi-square test statistic can be shown to be equal to 0.519, and the p -value is 0.471. Because the p -value is greater than 0.05 (the typical cutoff), you can't reject H_0 ; therefore you conclude that gender and opinion on the banning of cellphones while driving are independent and therefore not related. Your theory that gender had something to do with it just doesn't pan out; there's not sufficient evidence for it.



In general, *independence* means that you can find no major difference in the way the rows look as you move down a column. Put another way, the proportion of the data falling into each column across the row is about the same for each row. Because Table 14-4 has the same number of men as women, the row totals are the same, and you get the same expected cell counts for men and women in both the Agree column (10.5) and the Disagree column (49.5).

Comparing Two Tests for Comparing Two Proportions

You can use the Chi-square test to check whether two population proportions are equal. For example, is the proportion of female cellphone users the same as the proportion of male cellphone users?

You may be thinking, “But wait a minute, don’t statisticians already have a test for two proportions? I seem to remember it from my Stats I course . . . I’m thinking . . . yeah, it’s the Z-test for two proportions. What’s that test got to do with a Chi-square test?” In this section, you get an answer to that question and practice using both methods to investigate a possible gender gap in cellphone use.

Getting reacquainted with the Z-test for two population proportions

The way that most people figure out how to test the equality of two population proportions is to use a *Z-test for two population proportions*. With this test, you collect a random sample from each of the two populations, find and subtract their two sample proportions, and divide by their pooled standard error (see your Stats I textbook for details on this particular test).

This test is possible to do as long as the sample sizes from the two populations are large — at least five successes and five failures in each sample.

The null hypothesis for the Z-test for two population proportions is $H_0: p_1 = p_2$, where p_1 is the proportion of the first population that falls into the category of interest, and p_2 is the proportion of the second population that falls into the category of interest. And as always, the alternative hypothesis is one of the following choices, H_a : Not equal to, greater than, or less than.

Suppose you want to compare the proportion of male versus female cellphone users, where p_1 is the proportion of males who own a cellphone, and p_2 is the proportion of all

females who own a cellphone. You collect data, find the sample proportions from each group, take their difference and make a Z-statistic out of it using the formula

$$Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}, \text{ where } \hat{p} = \frac{x_1 + x_2}{n_1 + n_2}.$$

Here, x_1 and x_2 are the number of individuals from samples one and two, respectively, with the desired characteristic; n_1 and n_2 are the two sample sizes.

Suppose that you collect data on 100 men and 100 women and find 45 male cellphone owners and 55 female cellphone owners. This means that \hat{p}_1 equals $45 \div 100 = 0.45$, and \hat{p}_2 equals $55 \div 100 = 0.55$. Your samples have at least five successes (having the desired characteristic; in this case, cellphone ownership) and five failures (not having the desired characteristic, which is cellphone ownership). So you compute the Z-statistic for comparing the two population proportions (males versus females) based on this data; it's -1.41 , as shown on the last line of the Minitab output in Figure 14-3.

Figure 14-3:
Minitab
output
comparing
proportion of
male and
female
cellphone
owners.

Test Cellphone for Two Proportions

Sample	X	N	Sample p
M	45	100	0.450000
F	55	100	0.550000

Difference = p (1) - p (2)
Estimate for difference: -0.1
95% CI for difference: (-0.237896, 0.0378957)
Test for difference = 0 (vs not = 0): Z = -1.41 P-Value = 0.157

The p -value for the test statistic of $Z = -1.41$ is 0.157 (calculated by Minitab, or by looking at the area below the Z -value of -1.41 on a Z -table, which you should have in your Stats I text). This p -value (0.157) is greater than the typical α level (predetermined cutoff) of 0.05, so you can't reject H_0 . You can't say that the two population proportions aren't equal, so you must conclude that the proportion of male cellphone owners is no different than females.

Even though the sample seemed to have evidence for a difference (after all, 45 percent isn't equal to 55 percent), you don't have enough evidence in the data to say that this same difference carries over to the population. So you can't lay claim to a gender gap in cellphone use, at least not with this sample.

Equating Chi-square tests and Z-tests for a two-by-two table

Here's the key to relating the Z-test to a Chi-square test for independence. The Z-test for two proportions and the Chi-square test for independence in a two-by-two table (one with two rows and two columns) are equivalent if the sample sizes from the two populations are large enough — that is, when the number of successes and the number

of failures in each cell of the two samples is at least five.

If you use the Z-test to see whether the proportion of male cellphone owners is equal to the proportion of female cellphone owners, you're really looking at whether you can expect the same proportion of cellphone owners despite gender (after you take the sample sizes into account). And that means you're testing whether gender (male or female) is independent of cellphone ownership (yes or no).

If the proportion of female cellphone owners equals the proportion of male cellphone owners, the proportion of cellphone owners is the same regardless of gender, so gender and cellphone ownership are independent. On the other hand, if you find the proportion of male cellphone owners to be unequal to the proportion of female cellphone owners, you can say that cellphone use differs by gender, so gender and cellphone ownership are dependent.

With the cellphone data, you have 45 males using cellphones (out of 100 males) and 55 females using cellphones (out of 100 females). The Minitab output for the Chi-square test for independence (complete with observed and expected cell counts, degrees of freedom, test statistic, and p -value) is shown in Figure 14-4. The p -value for this test is 0.157, which is greater than the typical α level 0.05, so you can't reject H_0 .

Because the Chi-square test for independence and the Z-test tests are equivalent when you have a two-by-two table, the p -value from the Chi-square test for independence is identical to the p -value from the Z-test for two proportions. If you compare the p -values from Figures 14-3 and 14-4, you can see that for yourself.

Figure 14-4:
Minitab
output testing
independence
of gender
and
cellphone
ownership.

Chi-Square Test: Gender, Cellphone			
Expected counts are printed below observed counts			
Chi-Square contributions are printed below expected counts			
	Y	N	Total
M	45	55	100
	50.00	50.00	
	0.500	0.500	
F	55	45	100
	50.00	50.00	
	0.500	0.500	
Total	100	100	200
Chi-Sq = 2.000, DF = 1, P-Value = 0.157			

Also, note that if you take the Z-test statistic for this example (from Figure 14-3), which is -1.41 , and square it, you get 2.00, which is equal to the Chi-square test statistic for the same data (last line of Figure 14-4). It's also the case that the square of the Z-test statistic (when testing for the equality of two proportions) is equal to the corresponding Chi-square test statistic for independence.



The Chi-square test and Z-test are equivalent only if the table is a two-by-two table (two rows and two columns) and if the Z-test is two-tailed (the alternative hypothesis is that the two proportions aren't equal, instead of using Ha: One proportion is greater than or less than the other). If the Z-test isn't two-tailed, a Chi-square test isn't appropriate. If the two-way table has more than two rows or columns, use the Chi-square test for independence (because many categories mean you no longer have only two proportions, so the Z-test isn't applicable).

The car accident–cellphone connection

Researchers are doing a great deal of study of the effects of cellphone use while driving. One study published in the *New England Journal of Medicine* observed and recorded data in 1997 on 699 drivers who had cellphones and were involved in motor vehicle collisions resulting in substantial property damage but no personal injury. Each person's cellphone calls on the day of the collision and during the previous week were analyzed through the use of detailed billing records. A total of 26,798 cellphone calls were made during the 14-month study period.

One conclusion the researchers made was that "... the risk of a collision when using a cellphone is four times higher than the risk of a collision when a cellphone was not being used." They basically conducted a Chi-square test to see whether cellphone use and having a collision are independent, and when they found out the events were not, the researchers were able to examine the relationship further using appropriate ratios. In particular, they found that the risk of a collision is four times higher for those drivers using cellphones than for those who aren't.

Researchers also found out that the relative risk was similar for drivers who differed in personal characteristics, such as age and driving experience. (This finding means that they conducted similar tests to see whether the results were the same for drivers of different age groups and drivers of different levels of experience, and the results always came out about the same. Therefore, age and the experience of the driver weren't related to the collision outcome.)

The research also shows that "... calls made close to the time of the collision were found to be particularly hazardous ($p < 0.001$). Hands-free cellphones offered no safety advantage over hand-held units (p -value not significant)." **Note:** The items in parentheses show the typical way that researchers report their results: using p -values. The p in both cases of parentheses represents the p -value of each test.

In the first case, the p -value is very tiny, less than 0.001, indicating strong evidence for a relationship between collisions and cellphone use at the time. The second p -value in parentheses was stated to be insignificant, meaning that it was substantially more than 0.05, the usual α level. This second result indicates that using hands-free equipment didn't affect the chances of a collision happening; the proportion of collisions using hands-free cellphones versus using regular cellphones were found to be statistically the same (they could have easily occurred by chance under independence). Whether you use a regular or hands-free cellphone, may this study be a lesson to you!

Chapter 15

Using Chi-Square Tests for Goodness-of-Fit (Your Data, Not Your Jeans)

In This Chapter

- ▶ Understanding what goodness-of-fit really means
 - ▶ Using the Chi-square model to test for goodness-of-fit
 - ▶ Looking at the conditions for goodness-of-fit tests
-

Many phenomena in life may appear to be haphazard in the short term, but they actually occur according to some preconceived, preselected, or predestined model over the long term. For example, even though you don't know whether it will rain tomorrow, your local meteorologist can give you her model for the percentage of days that it rains, snows, is sunny, or cloudy, based on the last five years. Whether or not this model is still relevant this year is anyone's guess, but it's a model nonetheless. As another example, a biologist can produce a model for predicting the number of goslings raised by a pair of geese per year, even though you have no idea what the pair in your backyard will do. Is his model correct? Here's your chance to find out.

In this chapter, you build models for the proportion of outcomes that fall into each category for a categorical variable. You then test these models by collecting data and comparing what you observe in your data to what you expect from the model. You do this evaluation through a goodness-of-fit test that's based on the Chi-square distribution. In a way, a goodness-of-fit test is likened to a reality check of a model for categorical data.

Finding the Goodness-of-Fit Statistic

The general idea of a *goodness-of-fit* procedure involves determining what you expect to find and comparing it to what you actually observe in your own sample through the use of a test statistic. This test statistic is called the *goodness-of-fit test statistic* because it measures how well your model (what you expected) fits your actual data (what you observed).

In this section, you see how to figure out the numbers that you should expect in each category given your proposed model, and you also see how to put those expected values together with your observed values to form the goodness-of-fit test statistic.

What's observed versus what's expected

For an example of something that can be observed versus what's expected, look no

further than a bag of tasty M&M'S Milk Chocolate Candies. A ton of different kinds of M&M'S are out there, and each kind has its own variation of colors and tastes. For this study, any reference I give to M&M'S is to the original milk chocolate candy — my favorite.

The percentage of each color of M&M'S that appear in a bag is something Mars (the company that makes M&M'S) spends a lot of time thinking about. Mars wants specific percentages of each color in its M&M'S bags, which it determines through comprehensive marketing research based on what people like and want to see. Mars then posts its current percentages for each color of M&M'S on its Web site. Table 15-1 shows the percentage of M&M'S of each color in 2006.

Table 15-1 Expected Percentage of Each Color of M&M'S Milk Chocolate Candies (2006)

<i>Color</i>	<i>Percentage</i>
Brown	13%
Yellow	14%
Red	13%
Blue	24%
Orange	20%
Green	16%

Now that you know what to expect from a bag of M&M'S, the next question is, how does Mars deliver? If you were to open a bag of M&M'S right now, would you get the percentages of each color that you're supposed to get? You know from your previous studies in statistics that sample results vary (for a quick review of this idea, see Chapter 3). So you can't expect each bag of M&M'S to have exactly the correct number of each color of M&M'S as listed in Table 15-1. However, in order to keep customers happy, Mars should get close to the expectations. How can you determine how close the company does get?

Table 15-1 tells you what percentages are expected to fall into each category in the entire population of all M&M'S (that means every single M&M'S Milk Chocolate Candy that's currently being made). This set of percentages is called the *expected model* for the data. You want to see whether the percentages in the expected model are actually occurring in the packages you buy. To start this process, you can take a sample of M&M'S (after all, you can't check every single one in the population) and make a table showing what percentage of each color you observe. Then you can compare this table of observed percentages to the expected model.



Some expected percentages are known, as they are for the M&M'S, or you can figure them out by using math techniques. For example, if you're examining a single

die to determine whether or not it's a fair die, you know that if the die is fair, you should expect $\frac{1}{6}$ of the outcomes to fall into each category of 1, 2, 3, 4, 5, and 6.

As an example, I examined one 1.69-ounce bag of plain, milk-chocolate M&M'S (tough job, but someone had to do it), and you can see my results in Table 15-2, column two. (Think of this bag as a random sample of 56 M&M'S, even though it's not technically the same as reaching into a silo filled with M&M'S and pulling out a true random sample of 1.69 ounces. For the sake of argument, one bag is okay.)

Table 15-2 Percentage of M&M'S Observed in One Bag (1.69 oz.) Versus Percentage Expected

Color	Percentage Observed	Percentage Expected
Brown	$\frac{4}{56} = 7.14$	13.00
Yellow	$\frac{10}{56} = 17.86$	14.00
Red	$\frac{4}{56} = 7.14$	13.00
Blue	$\frac{10}{56} = 17.86$	24.00
Orange	$\frac{15}{56} = 26.79$	20.00
Green	$\frac{13}{56} = 23.21$	16.00
TOTAL	100.00	100.00

Compare what I observed in my sample (column two of Table 15-2) to what I expected to get (column three of Table 15-2). Notice that I observed a lower percentage of brown and red M&M'S than expected and a lower percentage of blues than expected. I also observed a higher percentage of yellow, orange, and green M&M'S than expected. Sample results vary by random chance, from sample to sample, and the difference I observed may just be due to this chance variation. But could the differences indicate that the expected percentages reported by Mars aren't being followed?

It stands to reason that if the differences between what you observed and what you expected are small, you should attribute that difference to chance and let the expected model stand. On the other hand, if the differences between what you observed and what you expected are large enough, you may have enough evidence to indicate that the expected model has some problems. How do you know which conclusion to make? The operative phrase is, "if the differences are large enough." You need to quantify this term *large enough*. Doing so takes a bit more machinery, which I cover in the next section.

Calculating the goodness-of-fit statistic

The goodness-of-fit statistic is one number that puts together the total amount of difference between what you expect in each cell compared to the number you observe. The term *cell* is used to express each individual category within a table format. For example, with the M&M'S example, the first columns of Tables 15-1 and 15-2 contain six cells, one for each color of M&M. For any cell, the number of items you observe in that cell is called the *observed cell count*. The number of items you expect in that cell (under the given model) is called the *expected cell count*. You get the expected cell count by taking the expected cell percentage times the sample size.

The expected cell count is just a proportion of the total, so it doesn't have to be a whole number. For example, if you roll a fair die 200 times, you should expect to roll ones $\frac{1}{6}$, or 16.67 percent, of the time. In terms of the number of ones you expect, it should be $0.1667 * 200 = 33.33$. Use the 33.33 in your calculations for goodness-of-fit; don't round to a whole number. Your final answer is more accurate that way.



The reason the goodness-of-fit statistic is based on the *number* in each cell rather than the *percentage* in each cell is because percents are a bit deceiving. If you know that 8 out of 10 people support a certain view, that's 80 percent. But 80 out of 100 is also 80 percent. Which one would you feel is a more-precise statistic? The 80 out of 100 percent because it uses more information. Using percents alone disregards the sample size. Using the counts (the number in each group) keeps track of the amount of precision you have.

For example, if you roll a fair die, you expect the percentage of ones to be $\frac{1}{6}$. If you roll that fair die 600 times, the expected number of ones will be $\frac{1}{6} * 600 = 100$. That number (100) is the expected cell count for the cell that represents the outcome of one. If you roll this die 600 times and get 95 ones, then 95 is the observed cell count for that cell.

The formula for the goodness-of-fit statistic is given by the following: $\sum_{\text{all cells}} \frac{(O-E)^2}{E}$, where E is the expected number in a cell and O is the observed number in a cell. The steps for this calculation are as follows:

1. **For the first cell, find the expected number for that cell (E) by taking the percentage expected in that cell times the sample size.**
2. **Take the observed value in the first cell (O) minus the number of items that are expected in that cell (E).**
3. **Square that difference.**
4. **Divide the answer by the number that's expected in that cell, (E).**
5. **Repeat steps one through four for each cell.**
6. **Add up the results to get the goodness-of-fit statistic.**



The reason you divide by the expected cell count in the goodness-of-fit statistic (step four) is to take into account the magnitude of any differences you find. For example, if you expected 100 items to fall in a certain cell and you got 95, the difference is 5. But in terms of a percentage, this difference is only $\frac{5}{100} = 5\%$ percent. However, if you expected 10 items to fall into that cell and you observed 5 items, the difference is still 5, but in terms of a percentage, it's $\frac{5}{10} = 50\%$ percent. This difference is

much larger in terms of its impact. The goodness-of-fit statistic operates much like a percentage difference. The only added element is to square the difference to make it positive. (That's done because whether you expected 10 and got 15 or expected 10 and got 5 makes no difference to others; you're still off by 50 percent.)

Table 15-3 shows the step-by-step calculation of the goodness-of-fit statistic for the M&M'S example, where O indicates observed cell counts and E indicates expected cell counts. To get the expected cell counts, you take the expected percentages shown in Table 15-1 and multiply by 56 because 56 is the number of M&M'S I had in my sample. The observed cell counts are the ones found in my sample, shown in Table 15-2.

Table 15-3 Goodness-of-Fit Statistic for M&M'S Example					
Color	O	E	$O - E$	$(O - E)^2$	$\frac{(O - E)^2}{E}$
Brown	4	$0.13 * 56 = 7.28$	$4 - 7.28 = -3.28$	10.76	1.48
Yellow	10	$0.14 * 56 = 7.84$	$10 - 7.84 = 2.16$	4.67	0.60
Red	4	$0.13 * 56 = 7.28$	$4 - 7.28 = -3.28$	10.76	1.48
Blue	10	$0.24 * 56 = 13.44$	$10 - 13.44 = -3.44$	11.83	0.88
Orange	15	$0.20 * 56 = 11.20$	$15 - 11.20 = 3.80$	14.44	1.29
Green	13	$0.16 * 56 = 8.96$	$13 - 8.96 = 4.04$	16.32	1.82
TOTAL	56			7.55	

The goodness-of-fit statistic for the M&M'S example turns out to be 7.55, the bolded number in the lower-right corner of Table 15-3. This number represents the total squared difference between what I expected and what I observed, adjusted for the magnitude of each expected cell count. The next question is how to interpret this value of 7.55. Is it large enough to indicate that colors of M&M'S in the bag aren't following the percentages posted by Mars? The next section addresses how to make sense of these results.

Interpreting the Goodness-of-Fit Statistic Using a Chi-Square

After you get your goodness-of-fit statistic, your next job is to interpret it. To do this, you need to figure out the possible values you could have gotten and where your statistic fits in among them. You can accomplish this task with a Chi-square goodness-of-fit test.

The values of a goodness-of-fit statistic actually follow a Chi-square distribution with $k - 1$ degrees of freedom, where k is the number of categories in your particular population (see Chapter 14 for the full details on the Chi-square). You use the Chi-square table (Table A-3 in the appendix) to find the p -value of your Chi-square test statistic.

If your Chi-square goodness-of-fit statistic is large enough, you conclude that the original model doesn't fit and you have to chuck it; there's too much of a difference between what you observed and what you expected under the model. However, if your goodness-of-fit

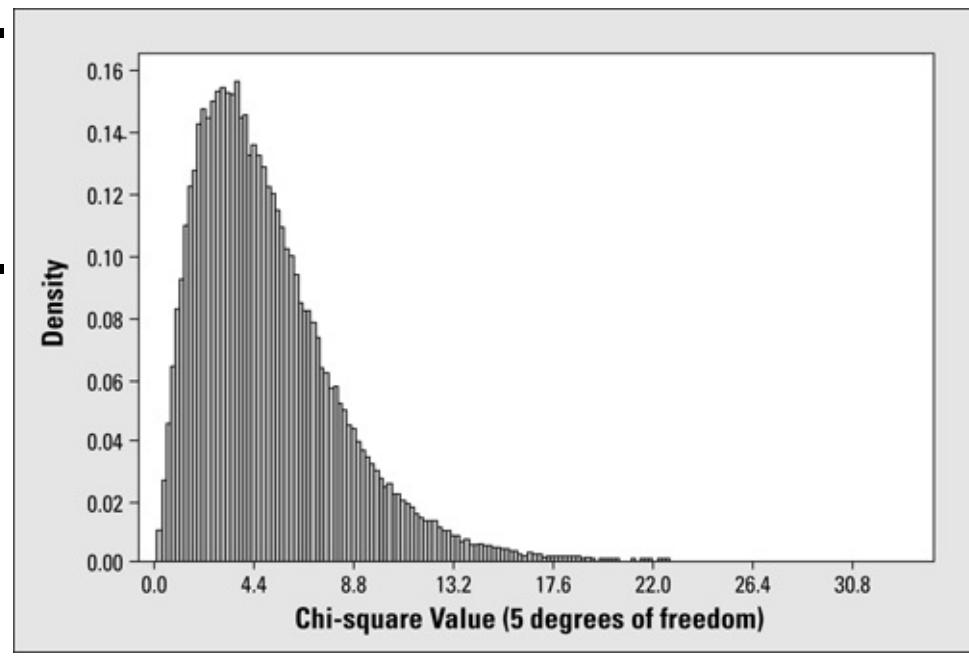
statistic is relatively small, you don't reject the model. (What constitutes a large or small value of a Chi-square test statistic depends on the degrees of freedom.)



The goodness-of-fit statistic follows the main characteristics of the Chi-square distribution. The smallest-possible value of the goodness-of-fit statistic is zero. Continuing the example from the previous section, if the M&M'S in my sample followed the exact percentages found in Table 15-1, the goodness-of-fit statistic would be zero. That's because the observed counts and the expected counts would be the same, so the values of the observed cell count minus the expected cell count would all be zero.

The largest-possible value of Chi-square isn't specified, although some values are more likely to occur than others. Each Chi-square distribution has its own set of likely values, as you can see in Figure 15-1. This figure shows a simulated Chi-square distribution with $6 - 1 = 5$ degrees of freedom (relevant to the M&M'S example). It basically gives a breakdown of all the possible values you could have for the goodness-of-fit statistic in this situation and how often they occur. You can see in Figure 15-1 that a Chi-square test statistic of 7.55 isn't unusually high, indicating that the model for M&M'S colors probably can't be rejected. However, more particulars are needed before you can formally make that conclusion.

Figure 15-1:
Chi-square distribution with 5 degrees of freedom.



Checking the conditions before you start

Every statistical technique seems to have a catch, and this case is no exception. In order to use the Chi-square distribution to interpret your goodness-of-fit statistic, you have to be sure you have enough information to work with in each cell. The stats gurus usually recommend that the expected count for each cell turns out to be greater than or equal to five. If it doesn't, one option is to combine categories to increase the numbers.

In the M&M'S example, the expected cell counts are all above seven (see Table 15-3), so the conditions are met. If this weren't the case, you should have taken a larger sample size, because you calculate the expected cell counts by taking the expected percentage in that cell times the sample size. If you increase the sample size, you increase the expected cell count. A higher sample size also increases your chances of detecting a real deviation from the model. This idea is related to the power of the test (see Chapter 3 for information on power).



After you collect your data, it's not right to go back and take a new and larger sample. It's best to set up the appropriate sample size ahead of time, and you can do this by determining what sample size you need to get the expected cell counts to be at least five. For example, if you roll a fair die, you expect $\frac{1}{6}$ of the outcomes to be ones. If you only take a sample of six rolls, you have an expected cell count of $\frac{1}{6} * 6 = 1$, which isn't enough. However, if you roll the die 30 times, your expected cell count is $\frac{1}{6} * 30 = 5$, which is just enough to meet the condition.

The steps of the Chi-square goodness-of-fit test

Assuming the necessary condition is met (see the previous section), you can get down to actually conducting a formal goodness-of-fit test.

The general version of the null hypothesis for the goodness-of-fit test is H_0 : The model holds for all categories; versus the alternative hypothesis H_a : The model doesn't hold for at least one category. Each situation will dictate what proportions should be listed in H_0 for each category. For example, if you're rolling a fair die, you have H_0 : Proportion of ones = $1s = \frac{1}{6}$; proportion of twos = $2s = \frac{1}{6}$; . . . ; proportion of sixes = $6s = \frac{1}{6}$.

Following are the general steps for the Chi-square goodness-of-fit test, with the M&M'S example illustrating how you can carry out each step:

1. Write down H_0 using the percentages that you expect in your model for each category.

Using a subscript to indicate the proportion (p) of M&M'S you expect to fall into each category (see Table 15-1), your null hypothesis is $H_0: p_{\text{brown}} = 0.13, p_{\text{yellow}} = 0.14, p_{\text{red}} = 0.13, p_{\text{blue}} = 0.24, p_{\text{orange}} = 0.20$, and $p_{\text{green}} = 0.16$. All these proportions must hold in order for the model to be upheld.

2. Write your H_a : This model doesn't hold for at least one of the percentages.

Your alternative hypothesis, H_a , in this case, would be: One (or more) of the probabilities given in H_0 isn't correct. In other words, you conclude that at least one of the colors of M&M'S has a different proportion than what's stated in the model.

3. Calculate the goodness-of-fit statistic using the steps in the earlier section

“Calculating the goodness-of-fit statistic.”

The goodness-of-fit statistic for M&M’S, from the earlier section, is 7.55. As a reminder, you take the observed number in each cell minus the expected number in that cell, square it, and divide by the expected number in that cell. Do that for every cell in the table and add up the results. For the M&M’S example, that total is equal to 7.55, the goodness-of-fit statistic.

4. Look up the Chi-square distribution with $k - 1$ degrees of freedom, where k is the number of categories you have.

You compare this statistic (7.55) to the Chi-square distribution with $6 - 1 = 5$ degrees of freedom (because you have $k = 6$ possible colors of M&M’S). (See Table A-3 in the appendix)

Looking at Figure 15-1 you can see that the value of 7.55 is nowhere near the high end of this distribution, so you likely don’t have enough evidence to reject the model provided by Mars for M&M’S colors.

5. Find the p -value of your goodness-of-fit statistic.

You use a Chi-square table to find the p -value of your test statistic (see Table A-3 in the appendix). (For more info on the Chi-square distribution, refer to Chapter 14.)



Because the Chi-square table can only list a certain number of results for each of the degrees of freedom, the exact p -value for your test statistic may fall between two p -values listed on the table.

To find the p -value for the test statistic in the M&M’S example (7.55), find the row for 5 degrees of freedom on the Chi-square table (Table A-3 in the appendix) and look at the numbers (the degrees of freedom is $k - 1 = 6 - 1 = 5$, where k is the number of categories). You see that the number 7.55 is less than the first value in the row (9.24), which has a p -value of 0.10. (Find the p -value by looking at the column heading above the number.) So the p -value for 7.55, which is the area to the right of 7.55 on Figure 15-1, must be greater than 0.10, because 7.55 is to the left of 9.24 on that Chi-square distribution.



Many computer programs exist (online or via a graphing calculator) that will find exact p -values for a Chi-square test, saving time and headaches when you have access to them (the technology, not the headaches). Using one such online p -value calculator, I found that the exact p -value for the goodness-of-fit test for the M&M’S example (test statistic 7.55 with 5 degrees of freedom for Chi-square) is 0.1828. To find online p -value calculators, simply type the name of the distribution and the word “ p -value” in an Internet search engine. For this example, search “Chi-square p -value.”

6. If your p -value is less than your predetermined cutoff (α), reject H_0 ; the model doesn’t hold. If your p -value is greater than α , you can’t reject the model.

A typical value of α is 0.05. Some data analysts may use a higher value (up to 0.10), and others may go lower (for example, 0.010). See Chapter 3 for more information on

choosing α and comparing your p -value to it.

Going again to the M&M'S example, the p -value, 0.18, is greater than 0.05, so you fail to reject H_0 . You can't say the model is wrong. So, Mars does appear to deliver on the percentages of M&M'S of each color as advertised. At least, you can't say it doesn't. (I'm sure Mars already knew that.)



Although some hypothesis tests are two-sided tests, the goodness-of-fit test is always a *right-tailed test*. You're only looking at the right tail of the Chi-square distribution when you're doing a goodness-of-fit test. That's because a small value of the goodness-of-fit statistic means that the observed data and the expected model don't differ much, so you stick with the model. If the value of the goodness-of-fit statistic is way out on the right tail of the Chi-square distribution, however, that's a different story. That situation means the difference between what you observed and what you expected is larger than what you should get by chance, and therefore, you have enough evidence to say the expected model is wrong.



You use the Chi-square goodness-of-fit test to check to see whether a specified model fits. A specified model is a model in which each possible value of the variable x is listed, along with its associated probability according to the model. For example, if you want to test whether three local hospitals take in the same percentage of emergency room patients, you test $H_0: p_1 = p_2 = p_3$, where each p represents the percentage of ER patients going to each hospital, respectively. In this case each p must equal 0.30 if the hospitals share the ER load equally.

Part V

Nonparametric Statistics: Rebels without a Distribution

The 5th Wave

By Rich Tennant



"Ted and I spent over 120 man-hours together analyzing the survey data, and here's what we discovered: Ted borrows pens and never returns them, he intentionally squeaks his chair to annoy me, and, evidently, I talk in my sleep."

In this part . . .

Suppose you're driving home and one of the streets is blocked. What do you do? You back up and find another way to get home. Nonparametric statistics is that alternative route you take if the regular parametric statistical methods aren't allowed. What's more, this alternate route actually turns out to be better when the regular route isn't available. In this part, you find out just how much better nonparametric statistics is using the sign test, the signed rank test, and many more.

Chapter 16

Going Nonparametric

In This Chapter

- ▶ Understanding the need for nonparametric techniques
 - ▶ Distinguishing regular methods from nonparametric methods
 - ▶ Laying the groundwork: The basics of nonparametric statistics
-

Many researchers do analyses involving hypothesis tests, confidence intervals, Chi-square tests, regression, and ANOVA. But nonparametric statistics doesn't seem to gain the same popularity as the other methods. It's more in the background — an unsung hero, if you will. However, nonparametric statistics is, in fact, a very important and very useful area of statistics because it gives you accurate results when other, more common methods fail.

In this chapter, you see the importance of nonparametric techniques and why they should have a prominent place in your data-analysis toolbox. You also discover some of the basic terms and techniques involved with nonparametric statistics.

Arguing for Nonparametric Statistics

Nonparametric statistics plays an important role in the world of data analysis in that it can save the day when you can't use other methods. The problem is that researchers often disregard, or don't even know about, nonparametric techniques and don't use them when they should. In that case, you never know what kind of results you get; what you do know is they could very well be wrong.

In the following sections, you see the advantages and the flexibility of using a nonparametric procedure. You also find out just how minimal the downside is, which makes it a win-win situation most of the time.

No need to fret if conditions aren't met

Many of the techniques that you typically use to analyze data, including many shown in this book, have one very strong condition on the data that must be met in order to use them: The populations from which your data are collected typically require a normal distribution. Methods requiring a certain type of distribution (such as a normal distribution) in order to use them are called *parametric* methods.

The following are ways to help you decide whether a population has a normal distribution, based on your sample:

- ✓ You can graph the data using a histogram, and see whether it appears to have a bell shape (a mound of data in the middle, trailing down on each side).



To make a histogram in Minitab, enter your data into a column. Go to Graph>Histogram, and click OK. Click on your variable in the left-hand box, and it appears in the Graph Variables box. Click OK, and check out your histogram.

- ✓ You can make a normal probability plot, which compares your data to that of a normal distribution, using an x - y graph (similar to the ones used when you graph a straight line). If the data do follow a normal distribution, your normal probability plot will show a straight line. If the data don't follow a normal distribution, the normal probability plot won't show a straight line; it may show a curve off to one side or the other, for example.



To make a normal probability plot in Minitab, enter your data in a column. Go to Graph>Probability Plot, and click OK. Click on your variable in the left-hand column, and it appears in the Graph Variables column. Click OK, and you see your normal probability plot.

When you find that the normal distribution condition is clearly not met, that's where nonparametric methods come in. *Nonparametric methods* are those data-analysis techniques that don't require the data to have a specific distribution. Nonparametric procedures may require one of the following two conditions (and these are only in certain situations):

- ✓ The data come from a symmetric distribution (which looks the same on each side when you cut it down the middle).
- ✓ The data from two populations come from the same type of distribution (they have the same general shape).

Note also that the normal distribution centers solely on the mean as its main statistic (for example, the Z-value for the hypothesis test for one population mean is calculated by taking the data value, subtracting the mean, and dividing by the standard deviation). So the condition that the population has a normal distribution automatically says you're working with the mean. However, many nonparametric procedures work with the *median*, which is a much more flexible statistic because it isn't affected by *outliers* (extreme values either above or below the mean) or *skewness* (a peak on one side and a long tail on the other side) as the mean is.

The median's in the spotlight for a change

Many times a particular statistics question at hand revolves around the center of a population —that is, the number that represents a typical value, or a central value, in the

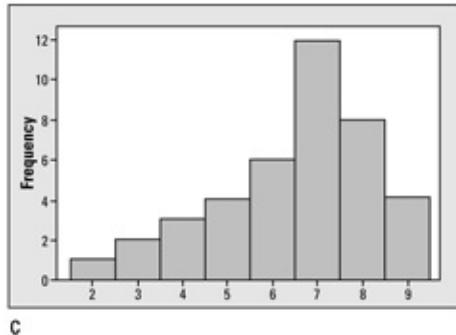
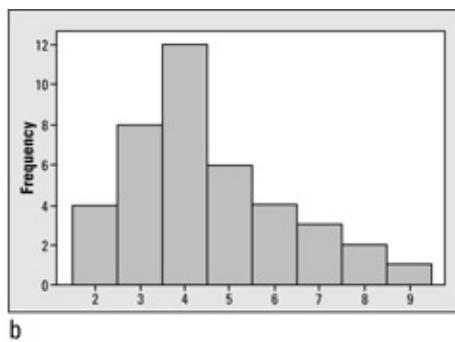
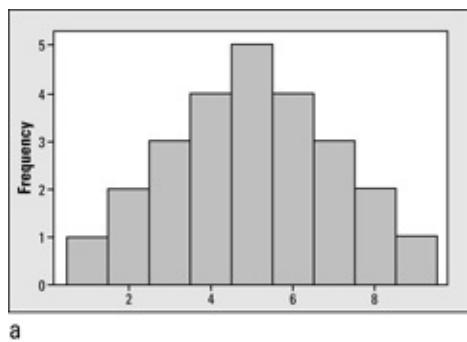
population. One of those measures of center is the *mean*. The *population mean* is the average value over the entire population, which is typically not known (that's why you take a sample). Many data analysts focus heavily on the population mean; they want to estimate it, test it, compare the means of two or more populations, or predict the mean value of a *y* variable given an *x* variable. However, the mean isn't the only measure of the center of a population; you also have the good ol' median.

You may recall that the *median* of a data set is the value that represents the exact middle when you order the data from smallest to largest. For example, in the data set 1, 5, 4, 2, 3, you order the data to get 1, 2, 3, 4, 5 and find that the number in the middle is 3, the median. If the data set has an even number of values, for example, 2, 4, 6, 8, then you average the two middle numbers to get your median — $(4 + 6) \div 2 = 5$ in this case.

As you may recall from Stats I, you can find the mean and the median of a data set and compare them to each other. You first organize your data into a histogram, and you look at its shape.

- ✓ **If the data set is symmetric**, meaning it looks the same on either side when you draw a line down the middle, the mean and median are the same (or close). Figure 16-1a shows an example of this situation. In this case, the mean and median are both 5.
- ✓ **If the histogram is skewed to the right**, meaning that you have a lot of smaller values and a few larger values, the mean increases due to those few large values, but the median isn't affected. In this case, the mean is larger than the median. Figure 16-1b shows an example of this situation, in which the mean is 4.5 and the median is 4.0.
- ✓ **If the histogram is skewed to the left**, you have many larger values that pile up, but only a few smaller values. The mean goes down because of the few small values, but the median still isn't affected. In this case, the mean is lower than the median. Figure 16-1c illustrates this situation with a 6.5 mean and a 7.0 median.

Figure 16-1:
Symmetric
and skewed
histograms.



My point is that the median is important! It's a measure of the center of a population or a sample data set. The median competes with the mean and often wins. Researchers use nonparametric procedures when they want to estimate, test, or compare the median(s) of one or more populations. They also use the median in cases where their data are symmetric but don't necessarily follow a normal distribution, or when they want to focus on a measure of center that's not influenced by outliers or skewness.

For example, if you look at house prices in your neighborhood, you may find a large number of houses within a certain relatively small price range as well as a few homes that cost a great deal more. If a real estate agent wants to sell a house in your neighborhood and intends to justify a high price for it, she may report the mean price of homes in your neighborhood because the mean is affected by outliers. The mean is higher than the median in this case. But if the agent wants to help someone buy a house in your 'hood, she looks at the median of the house prices in the neighborhood because the median isn't affected by those few higher-priced homes and is lower than the mean.

Now suppose you want to come up with a number that describes the typical house price in your entire county. Should you use the mean or the median? You gathered techniques in Stats I for estimating the mean of a population (see Chapter 3 for a quick review), but you probably didn't hear about how to come up with a confidence interval for the median of a population. Oh sure, you can take a random sample and calculate the median of that sample. But you need a margin of error to go with it. And I'll tell you something — the formula for the margin of error for the mean doesn't work for the margin of error associated with the median. (Hang on, this book has you covered on that score.)

So, what's the catch?

You may be wondering, what's the catch if I use a nonparametric technique? A downside must be lurking around here somewhere. Well, many researchers believe that nonparametric techniques water down statistical results; for example, suppose you find an actual difference between two population means, and the populations really do have a normal distribution. A parametric technique, the hypothesis test for two means, would likely detect this difference (if the sample size was large enough).

The question is, if you use a nonparametric technique (which doesn't need the populations to be normal), do you risk the chance of not finding the difference? The answer is maybe, but the risk isn't as big as you think. More often than not, nonparametric procedures are only slightly less efficient than parametric procedures (meaning they don't work quite as well at detecting a significant result or at estimating a value as parametric procedures) when the normality condition is met, but this difference in efficiency is small.

But the big payoff occurs when the normal distribution conditions aren't met. Parametric techniques can make the wrong conclusion, and corresponding nonparametric techniques can lead to a correct answer. Many researchers don't know this, so spread the word!



The bottom line: Always check for normality first. If you're very confident that the normality condition is met, go ahead and use parametric procedures because they're more precise. If you have any doubt about the normality condition, use nonparametric procedures. Even if the normality condition is met, nonparametric procedures are only a little less precise than parametric procedures. If the normality condition isn't met, nonparametrics provide appropriate and justifiable results where parametric procedures may not.

Mastering the Basics of Nonparametric Statistics

Because you may not have run into nonparametric statistics during your Stats I class, your first step toward using these techniques is figuring out some of the basics. In this section, you get to know some of the terminology and major concepts involved in nonparametric statistics. These terms and concepts are commonly used in Chapters 17 through 20 of this book.

Sign

The *sign* is a value of 0 or 1 that's assigned to each number in the data set. The sign for a value in the data set represents whether that data value is larger or smaller than some specified number. The value of +1 is given if the data value is greater than the specified number, and the value of 0 is given if the data value is less than or equal to the specified number. For example, suppose your data set is 10, 12, 13, 15, 20, and your specified number for comparison is 16. Because 10, 12, 13, and 15 are all less than 16, they each receive a sign of 0. Because 20 is greater than 16, it receives a sign of +1.

Several uses of the sign statistic appear in nonparametric statistics. You can use signs to test to see if the median of a population equals some specified value. Or you can use signs to analyze data from a matched-pairs experiment (where subjects are matched up according to some variable and a treatment is applied and compared). You also can use signs in combination with other nonparametric statistics. For example, you can combine signs with ranks to develop statistics for comparing the median of two populations. (Ranks are discussed in the next section and are used in a hypothesis test for two population medians in Chapter 18.)

In the following sections, you see exactly how to use the sign statistic to test the median of a population and analyze data in a matched-pairs experiment.

Testing the median

You can use signs to test whether the median of a population is equal to some value m . You do this by conducting a hypothesis test based on signs. You have $H_0: \text{Median} = m$ versus $H_a: \text{Median} \neq m$ (or, you can use a $>$ or $<$ sign in H_a also). Your test statistic is the sum of the signs for all the data. If this sum is significantly greater or significantly smaller than what's expected if H_0 were true, you reject H_0 . Exactly how large or how small the sum of the signs must be to reject H_0 is given by the sign test (refer to Chapter 17).

Suppose you're testing whether the median of a population is equal to 5. That is, you're testing $H_0: \text{Median} = 5$ versus $H_a: \text{Median} \neq 5$. You collect the following data: 4, 4, 3, 3, 2, 6, 4, 3, 3, 5, 7, 5. Ordering the data, you get 2, 3, 3, 3, 3, 4, 4, 4, 5, 5, 6, 7. Now you find the sign for each value in the data set, determined by whether the value is greater than 5. The sign of the first data value, 2, is 0, because it's below 5. Each of the 3s receives a sign of 0, as do the 4s and 5s, for the same reason. Only the numbers 6 and 7 receive a sign of +1, being the only values in the data set that are greater than 5 (the number of interest for the median).

By summing the signs, you're in essence counting the number of values in the data set that are greater than the given quantity in H_0 . For example, the total of all the signs of the ordered data values is

$$0 + 0 + 0 + 0 + 0 + 0 + 0 + 0 + 1 + 1 = 2$$

You can see that the total number of data values above 5 (the number of interest for the median) is 2. The fact that the total of the signs (2) is much less than half the sample size

gives you some evidence that the median is probably not 5 here because the median represents the middle of the population. If the median were truly 5 in the population, your sample should yield about six values below it and six values above.

Doing a matched-pairs experiment

You can use signs in a *matched-pairs experiment*, where you use the same subject twice or pair up subjects on some important variables. For example, you can use signs to test whether or not a certain treatment resulted in an improvement in patients, compared to a control. In the cases where the sign statistic is used, improvement is measured not by the mean of the differences in the responses for treatment versus control (as in a paired *t*-test), but rather by the median of the differences in the responses.

Suppose you're testing a new antihistamine for allergy patients. You take a sample of 100 patients and have each patient assess the severity of his or her allergy symptoms before and after taking the medication on a scale from 1 (best) to 10 (worst). (Of course, you do a controlled experiment in which some of the patients get a placebo to adjust for the fact that some people may perceive their symptoms to be going away just because they took something.)

In this study, you're not interested in what level the patients' symptoms are at, but rather in how many patients had a lower level of symptoms after taking the medicine. So you take the symptom level before the experiment minus the symptom level after the experiment for each subject.

- ✓ If that difference is positive, the medicine appears to have helped, and you give that person a sign of +1 (in other words, count them as a success).
- ✓ If the difference is zero, the medicine had no effect, and you give that person a sign of 0.

Remember, though, that the difference could be negative, indicating that the symptoms before were lower than the symptoms after; in other words, the medicine made their symptoms worse. This scenario results in a sign of 0 as well.

After you've found the sign for each value or pair in the data set, you're ready to analyze it by using the sign test or the signed rank test (see Chapter 17).

Rank

Ranks are a nice way to use important information from a data set without using the actual values of the data themselves. Rank comes into play in nonparametric statistics when you're not interested in the values of the data, but rather where they stand compared to some supposed value for the median or compared to the ranks of values in another data set from another population. (You can see ranks in action in Chapter 18.)

The *rank* of a value in a data set is the number that represents its place in the ordering,

from smallest to largest, within the data set. For example, if your data set is 1, 10, 4, 2, 1,000, you can assign the ranks in the following way: 1 gets the rank of 1 (because it's the smallest), 2 gets the rank of 2, 4 gets the rank of 3 (being the third-smallest number in the ordered data set), 10 gets the rank of 4, and 1,000 gets the rank of 5 (being the largest).

Now suppose your data set is 1, 2, 20, 20, 1,000. How do you assign the ranks? You know that 1 gets the rank of 1 (being the smallest), 2 gets the rank of 2, and 1,000 gets the rank of 5 (being the largest). But what about the two 20s in this data set? Should the first 20 get a rank of 3 and the second 20 get the rank of 4? That order doesn't seem to make sense, because you can't distinguish between the two 20s.



When two values in a data set are the same, it's called a tie. To assign ranks when there are ties, you take the average of the two ranks that the values need to fill and assign each tied value that average rank. If you have a tie between three numbers, you have three ranks, so take the sum of the ranks divided by three.

In this case, because both 20s are vying for the ranks of 3 and 4, assign each of them the rank of 3.5, the average of the two ranks they must share. I show the final ranking for the data set 1, 2, 20, 20, 1,000 in Table 16-1.

Table 16-1 Ranks of the Values in the Data Set 1, 2, 20, 20, 1,000

Data Value	Rank Assigned
1	1
2	2
20	3.5
20	3.5
1,000	5



The lowest a rank can be is 1, and the highest a rank can be is n , where n is the number of values in the data set. If you have a negative value in a data set, for example, if your data set is $-1, -2, -3$, you still assign the ranks 1 through 3 to those data values. Never assign negative ranks to negative data. (By the way, when you order the data set $-1, -2, -3$, you get $-3, -2, -1$, so -3 gets the rank of 1, -2 gets the rank of 2, and -1 gets the rank of 3.)

Signed rank

A *signed rank* combines the idea of the sign and the rank of a value in a data set, with a small twist. The sign indicates whether that number is greater than, less than, or equal to

a specified value. The rank indicates where that number falls in the ordering of the data set from smallest to largest.

To calculate the signed rank for each value in the data set, follow these steps:

1. Assign a sign of +1 or 0 to each value in the data set, according to whether it's greater than some value specified in the problem.

If it's greater than the specified value, give it a sign of +1; if it's less than or equal to the specified value, give it a sign of 0.

2. Rank the original data from smallest to largest, according to their absolute values.



Statisticians call these values the *absolute ranks*. The absolute value of any number is the positive version of that number. The notation for absolute value is $| |$, with the number between those lines. For example, $|-2| = 2$, $|+2| = 2$, and $|0| = 0$.

3. Multiply the sign times the absolute rank to get the signed rank for each value in the data set.

One scenario in which you can use signed ranks is an experiment in which you compare a response variable for a treatment group versus a control group. You can test for difference due to a treatment by collecting the data in pairs, either both from the same person (pretest versus post-test) or from two individuals who are matched up to be as similar as possible.

For example, suppose you compare four patients regarding their weight loss on a diet program. You're really wondering whether the overall change in weight is less than zero for the population. The following two factors are important:

- ✓ Whether or not the person lost weight
- ✓ How the person's weight change measures up, compared to everyone else in the data set

You measure the person's weight before the program (the pretest) as well as his weight after the program (the post-test). The change is the important facet of the data you're interested in, so you apply the signs to the changes in weight. You give the change a sign of +1 if the person lost weight (constituting a success for the program) and a sign of 0 if the person stayed the same or gained weight (thus not contributing to the success of the program). You convert all the changes in weight loss to their absolute values, and then you rank the absolute values (in other words, you've found the absolute ranks of the changes in weight). The signed rank is the product of the sign and the absolute rank. After determining the signed rank, you can really compare the effectiveness of the program: Large signed ranks indicate a big weight loss.

For example, weight changes of -20, -10, +1, and +5 have signs of +1, +1, 0, and 0. The

absolute values of the weight changes are 20, 10, 1, and 5. Their absolute ranks, respectively, are 4, 3, 1, and 2. The signed ranks are $4 * 1 = 4$, $3 * 1 = 3$, $1 * 0 = 0$, and $2 * 0 = 0$.

Rank sum

A *rank sum* is just what it sounds like: The sum of all the ranks. You typically use rank sums in situations when you're comparing two or more populations to see whether one has a central location that's higher than the other. (In other words, if you looked at the populations in terms of their histograms, one would be shifted to the right of the other on the number line.)

Here's a way in which researchers use rank sums: Suppose you're comparing quiz scores for two classes, and they don't have a normal distribution, hence you want to use nonparametric techniques to compare them. The total possible points on this quiz is 30. You collect random samples of five quiz scores from each of the classes. Suppose you collect the following sample data:

Class Number One	Class Number Two
22	23
23	30
20	27
25	28
26	25

The twist here is to combine all the data into one big data set, rank all the values, and sum the ranks for the first sample and then the second sample. Then compare the two rank sums. If one rank sum is higher, this outcome may indicate that a particular class did better on the quiz.

For the quiz example, the ordered data for the combined classes appear in the first line, with the respective ranks appearing in the second line. Circled scores came from the first class.

Ordered Data	(20)	(22)	(23)	23	(25)	25	(26)	27	28	30
Respective Ranks	1	2	3.5	3.5	5.5	5.5	7	8	9	10

The rank sum for the first class is $1 + 2 + 3.5 + 5.5 + 7 = 19$, which is quite a bit lower than the rank sum for the second class, $3.5 + 5.5 + 8 + 9 + 10 = 36$. This result tells you that the second class did better on the quiz than the first class, for this sample.

Chapter 18 shows you how to use a rank sum test to see whether the shapes of two population distributions are the same, meaning the values they take on and how often those values occur in each population. In Chapter 19, you can find even more uses for

rank sums, including the Kruskal-Wallis test.



Note that taking the mean of each data set and comparing them by using a two-sample t-test would be wrong in the quiz example because the quiz scores admittedly don't have a normal distribution. Indeed if the quiz were easy, you'd get many high scores and few low ones, and the population would be skewed left. On the other hand, if the quiz were hard, you'd get many low scores and few high ones, and the population would be skewed right (don't think too much about that scenario). In either case, you need a nonparametric procedure. See Chapter 18 for more on the nonparametric equivalent of the t-test.

Chapter 17

All Signs Point to the Sign Test and Signed Rank Test

In This Chapter

- ▶ Testing and estimating the median: The sign test
 - ▶ Figuring out when and how to use the signed rank test
-

The hypothesis tests you see in Stats I use well-known distributions like the normal distribution or the *t*-distribution (see Chapter 3). Using these tests requires that certain conditions be met, such as the type of data you’re using, the distribution of the population the data came from, and the size of your data set. Such procedures that involve such conditions are called *parametric procedures*. In general, parametric procedures are very powerful and precise, and statisticians use them as often as they can.

But, situations do arise in which your data don’t meet the conditions for a parametric procedure. Perhaps you just don’t have enough data (the biggest hurdle is whether the data come from a population with a normal distribution), or your data are just of a different type than quantitative data, such as ranks (where you don’t collect numerical data, but instead just order the data from low to high or vice versa).

In these situations, your best bet is a *nonparametric procedure* (see Chapter 16 for background info). In general, nonparametric procedures aren’t as powerful as parametric procedures, but they have very few assumptions tied to them. Moreover, nonparametric procedures are easy to carry out, and their formulas make sense. Most importantly, nonparametric procedures give accurate results compared to the use of parametric procedures when the conditions of parametric procedures aren’t met or aren’t appropriate.

In this chapter, you use the sign test and the Wilcoxon signed rank test to test or estimate the median of one population. These nonparametric procedures are the counterparts to the one-sample and matched pairs *t*-tests, which require data from a normal population.

Reading the Signs: The Sign Test

You use the one-sample *t*-test from Stats I to test whether or not the population mean is equal to a certain value. It requires the data to have a normal distribution. When this condition isn’t met, the *sign test* is a nonparametric alternative for the one-sample *t*-test. It tests whether or not the population median is equal to a certain value.

What makes the sign test so nice is that it's based on a very basic distribution, the binomial. You use the binomial distribution when you have a sequence of n trials of an experiment, with only two possible outcomes each time (success or failure). The probability of success is denoted by p , and is the same for each trial. The variable is x , the number of success in the n trials. (For more info on the binomial see your Stats I text.)

The only condition of the sign test is that the data are ordinal or quantitative — not categorical. However, this is no big deal because if you're interested in the median, you don't collect categorical data anyway.

Here are the steps for conducting the sign test. Note that Minitab can do steps four through eight for you; however, understanding what Minitab does behind the scenes is important, as always.

1. Set up your null hypothesis: $H_0: m = m_o$.

The true value of the median is m , and m_o is the claimed value of the median (the value you're testing).

2. Set up your alternative hypothesis. Your choices are $H_a: m \neq m_o$; or $H_a: m > m_o$; or $H_a: m < m_o$.

Which H_a you choose depends on what conclusion you want to make in the case that H_0 is rejected. For example, if you only want to know when the median is greater than some number m , use $H_a: m > m_o$. Chapter 3 tells you more about setting up alternative hypotheses.

3. Collect a random sample of (ordinal or quantitative) data from the population.

4. Assign a plus or minus sign to each value in the data set.

If an observation is less than m_o , assign it a minus (-) sign. If the observation is greater than m_o , give it a plus (+) sign. If the observation equals m_o , disregard it and let the sample size decrease by one.

In terms of the binomial distribution, you have n values in the data set, and each one has one of two outcomes: It falls either below m_o or above it. (This is akin to success and failure.)

5. Count up all the plus signs. This sum is your test statistic, noted by k .

In terms of the binomial, this sum represents the total number of successes, where a plus (+) sign is the designated success.

6. Locate the test statistic k (from step five) on the binomial distribution (using Table A-2 in the appendix).

You determine where your test statistic falls on the binomial distribution by looking it up in a binomial distribution table (check your textbook). To do this, you need to know n , k , and p .

Your sample size is n , your test statistic is k from step five, but what's your value of p , the probability of success? If the null hypothesis H_0 is true, 50 percent of the data should lie below m_o and 50 percent should lie above it. This corresponds to a success

(+) having a probability of $p = 0.50$ on the binomial distribution.

7. Find the *p*-value of your test statistic:

- If H_a has a $<$ sign, add up all the probabilities on the binomial table for $x \leq k$.
- If H_a has a $>$ sign, add up all the probabilities on the binomial table for $x \geq k$.
- If H_a has a \neq sign, add up the probabilities on the binomial table of x being greater than or equal to k and double this value. This gives you the *p*-value of the test.

8. Make your conclusion.

If the *p*-value from step six is less than the predetermined value of α (typically 0.05), reject H_0 and say the median is greater than, less than, or $\neq m_o$, depending on H_a . Otherwise, you can't reject H_0 .



To run a sign test in Minitab, enter your data in a single column. Go to Stat>Nonparametric>One-sample Sign. Click on your variable in the left-hand box, and click Select. The variable will appear in the Variables box. Then click OK, and you get the results of the sign test.

In the sections that follow, I show you two different ways in which you can use the sign test:

- ✓ To test or estimate the median of one population
- ✓ To test or estimate the median difference of data where the observations come in pairs, either from the same individual (pretest versus post-test) or individuals paired up according to relevant characteristics

Testing the median

Situations arise in which you aren't interested in the mean, but rather the median of a population. (Chapter 16 has more on the median.) For example, perhaps the data don't have a normal, or even a symmetric, distribution. When you want to estimate or test the median of a population (call it m), the sign test is a great option.

Suppose you're a real estate agent selling homes in a particular neighborhood, and you hear from other agents that the median house price in that neighborhood is \$110,000. You think the median is actually higher. Because you're interested in the median price of a home rather than the mean price, you decide to test the claim by using a sign test. Follow the steps of the sign test:

1. Set up your null hypothesis. Because the original claim is that the median price of a home is \$110,000, you have $H_0: m = \$110,000$.
2. Set up the alternative hypothesis. Because you believe the median is higher than \$110,000, your alternative hypothesis is $H_a: m > \$110,000$.
3. Take a random sample of ten homes in the neighborhood. You can see the data in

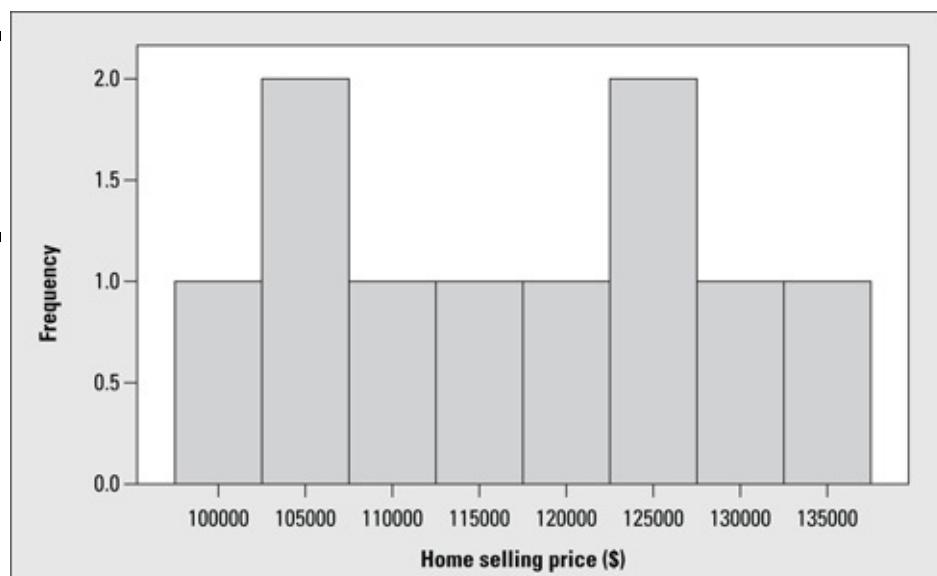
Table 17-1; its histogram is shown in Figure 17-1.

Now the question is, is the median selling price of all homes in the neighborhood equal to \$110,000, or is it more than that (as you suspect)?

Table 17-1 Sample of House Prices in a Neighborhood

House	Price	Sign (Compared to \$110,000)
1	\$132,000	+
2	\$107,000	-
3	\$111,000	+
4	\$105,000	-
5	\$100,000	-
6	\$113,000	+
7	\$135,000	+
8	\$120,000	+
9	\$125,000	+
10	\$126,000	+

Figure 17-1:
Histogram of
the selling
prices of ten
houses.



4. Assign a plus sign to any house price more than \$110,000 and a minus sign to any house less than \$110,000. (See column three of Table 17-1.)
5. Find your test statistic. Your test statistic is 7, the number of “+” signs in your data set (see Table 17-1), representing the number of houses in your sample whose prices are above \$110,000.
6. Compare your test statistic to the binomial distribution (refer to a binomial distribution table) to find the p -value.

For this case, look at the row in the binomial table where $n = 10$ (the sample size) and $k = 7$ (the test statistic) and the column where $p = 0.50$ (because if the population median equals m_o , 50 percent of the values in the population should be above it and 50 percent below it). According to the table, you find the probability that x equals 7 is 0.117.

Because you have a right-tailed test (meaning H_a has a $>$ sign in it), you add up the probabilities of being at or beyond 7 to get the p -value. The p -value in this case is $0.117 + 0.044 + 0.010 + 0.001 = 0.172$.

7. To conclude, compare the p -value (0.172) to the predetermined α (I always use 0.05). Because the p -value is greater than 0.05, you can't reject H_0 . You don't have enough evidence to say the median house selling price is more than \$110,000.

Figure 17-2 shows these results as calculated by Minitab.

Figure 17-2:
Sign test for
house prices
conducted
by Minitab.

Sign Test for Median: Selling Price						
Sign test of median = 110000 versus > 110000						
Selling	N	Below	Equal	Above	P	Median
Price	10	3	0	7	0.1719	116500



If your data are close to normal and the mean is the more appropriate measure of center for your situation, don't use the sign test. Instead, use the one-sample t -test (or Z-test). The sign test isn't quite as powerful (able to reject H_0 when it should) as the t -test in situations where the conditions for the t -test are met. More importantly, though, don't run to the t -test to reanalyze your data if the sign test doesn't reject H_0 . That would be improper and unethical. In general, statisticians consider the idea of following a nonparametric procedure with a parametric procedure in hopes of getting more significant results to be *data fishing*, which is analyzing data in different ways until a statistically significant result appears.

Estimating the median

You can also use the sign test to find a confidence interval for one population median. This comes in handy when you're interested in estimating what the median value of a population is, such as the median income of a household in the United States or the median salary of people fresh out of an MBA program.

Following are the steps for conducting a confidence interval for the median by using the test statistic for the sign test, assuming your random sample of data has already been collected. Note that Minitab can calculate the confidence interval for you (steps two to five), but knowing how Minitab does the steps is important:

1. **Determine your level of confidence, $1 - \alpha$ (that is, how confident you want to be that this process will correctly estimate m over the long term).**

The typical confidence level data analysts use is 95 percent (see Chapter 3 for more information).

2. **On the binomial table (Table A-2 in the appendix), find the section for n equal to your sample size, and the column where $p = 0.50$ (because the median is the point**

where 50 percent of the data lies below and 50 percent lies above).

You'll find probabilities for values of x from 0 to n in that section.

3. Starting at each end ($x = 0$ and $x = n$) and moving one step at a time toward the middle of the x values, add up the probabilities for those values of x until you pass the total of α (which is one minus your confidence level).

4. Record the number of steps that you had to make just before you passed the value of $1 - \alpha$. Call this number c .

5. Order your data set from smallest to largest. Starting at each end, work your way to the middle until you reach the c th number from the bottom and the c th number from the top.

6. Use these numbers as the low end and the high end of an interval. This result is your confidence interval for the median.

You can use these steps to find a confidence interval for the median in the house-price example from the preceding section. Here's how this example breaks down:

1. Let your confidence level be set at $1 - \alpha = 0.95$.

2. On the binomial table (Table A-2 in the appendix), look at the section where $n = 10$ (the sample size) and $p = 0.50$. These values are listed in Table 17-2.

Table 17-2 Binomial Probabilities to Help Calculate a Confidence Interval for the Median ($n = 10, p = 0.50$)

x	$p(x)$
0	0.001
1	0.010
2	0.044
3	0.117
4	0.205
5	0.246
6	0.205
7	0.117
8	0.044
9	0.010
10	0.001

3. Start with the outermost values of x ($x = 0$ and $x = 10$) and sum those probabilities to get $0.001 + 0.001 = 0.002$. Because you haven't yet passed 0.05 (the value of α), you go to the second-innermost values of x ($x = 1$ and $x = 9$). Add their probabilities to what you have so far to get 0.002 (old total) + $0.010 + 0.010 = 0.022$. You're still not past 0.05 (α), so go one more step. Add the third-innermost probabilities for $x = 2$ and $x = 8$ to the grand total to get 0.022 (old total) + $0.044 + 0.044 = 0.110$. You've now passed the value of $\alpha = 0.05$. The value of c equals 2 because you passed 0.05 at the third-innermost values of x ,

and you back off one step from there to get your value of c .

4. Order your data (Table 17-1) from smallest to largest, giving you (in dollars): 100,000, 105,000, 107,000, 111,000, 113,000, 120,000, 125,000, 126,000, 132,000, and 135,000.
5. Work your way in from each end of the data set to take the second-innermost values (because $c = 2$): the numbers \$105,000 and \$132,000. Put these two numbers together to form an interval, and you conclude that a 95 percent confidence interval for the median selling price for a home in this neighborhood is between \$105,000 and \$132,000.



To find a $1 - \alpha$ percent confidence interval for the median using Minitab based on the sign test, enter your data into a single column. Go to Stat>Nonparametrics>One-sample Sign. Click on the variable in the left-hand column for which you want the confidence interval, and it appears in the Variables column. Click the circle that says Confidence Interval, and type in the value of $1 - \alpha$ you want for your confidence level. (The default is 95 percent, written as 95.) Click OK to get the confidence interval.

Testing matched pairs

The most useful application of the sign test is in testing matched pairs of data — that is, data that come in pairs and represent two observations from the same person (pretests versus post-tests, for instance) or one set of data from each pair of people who are matched according to relevant characteristics. In this section, you see how you can compare data from a matched-pairs study to look for a treatment effect, using a sign test for the median.



The idea of using a sign test for the median difference with matched-pairs data is similar to using a t -test for the mean differences with matched-pairs data. (For details on matched-pairs data and the t -test, see your Stats I text.) You use a test of the median (rather than the mean) when the data don't necessarily have a normal distribution, or if you're only interested in the median difference rather than the mean difference.

First, you set up your hypothesis, H_0 : The median is zero (indicating no difference between the pairs). Your alternative hypothesis is H_a : The median is $\neq 0$, > 0 , or < 0 , depending on whether you want to know if the treatment made any difference, made a positive difference, or made a negative difference compared to the control. Then you collect your data (two observations per person or a pair of observations from each pair of people you've matched up). After that, you use Minitab to conduct steps four to seven of the sign test.

For example, suppose you wonder whether taking a test while chewing gum decreases

test anxiety. You pair 20 students according to relevant factors such as GPA, score on previous exams, and so on. One member of each pair is randomly selected to chew gum during the exam, and the other member of the pair doesn't. You measure test anxiety of each person via a very short survey right after they turn in their exams. You measure the results on a scale of 1 (lowest anxiety level) to 10 (highest anxiety level). Table 17-3 shows the data based on a sample of ten pairs.

Table 17-3 Testing the Effectiveness of Chewing Gum in Lowering Test Anxiety

Pair	Anxiety Level — Gum	Anxiety Level — No Gum	Difference (Gum/No Gum)	Sign
1	9	10	-1	-
2	6	8	-2	-
3	3	1	+2	+
4	3	5	-2	-
5	4	4	0	none
6	2	7	-5	-
7	2	6	-4	-
8	8	10	-2	-
9	6	8	-2	-
10	1	3	-2	-



The actual levels of test anxiety aren't important here; what matters is the difference between anxiety levels within each pair. So, instead of looking at all the individual anxiety levels, you can look at the difference in anxiety levels for each pair. This method gives you one data set, not two. (In this case, to calculate the differences in each pair, you can use the formula test anxiety without gum minus test anxiety with gum, and look for an overall difference that's positive.) Typically, in the case of matched-pairs data, you're testing whether the median difference equals zero. In other words, $H_0: m = 0$; the same holds in the test anxiety example.

The differences in anxiety levels for each pair in your data set now become a single data set (see column four of Table 17-3). You can now use the regular sign test methods to analyze this data, using $H_0: m = 0$ (no median difference in test anxiety of gum versus no gum) versus $H_a: m < 0$ (chewing gum reduces test anxiety).

Assign each difference a plus or minus sign, depending on whether it's greater than zero (plus sign) or less than zero (minus sign). Your test statistic is the total number of plus signs, 1, and the relevant sample size is $10 - 1 = 9$. (You don't count the data that hit the median of zero right on the head.)

Now compare this test statistic to the binomial distribution with $p = 0.50$ and $n = 9$, using the binomial table (Table A-2 in the appendix). You have a test statistic of $k = 1$, and you want to find the probability that $x \leq 1$ (because you have a left-tailed test, see step six of

the sign test from the earlier section “Reading the Signs: The Sign Test”). Under the column for $p = 0.50$ in the section for $n = 9$, you get the probability of 0.018 for $x = 1$ and 0.002 for $x = 0$. Add these values to get 0.020, your p -value. This result means that you reject H_0 at the predetermined α level of 0.05. This tells you the anxiety levels for gum versus no gum are different. Now, how are they different? Based on this data, you conclude that chewing gum on an exam appears to decrease test anxiety because there are more negative differences than positive differences.

Going a Step Further with the Signed Rank Test

The signed rank test is more powerful than the sign test at detecting real differences in the median. The most common use of the signed rank test is testing matched-pairs data for a median difference due to some treatment (like chewing-gum use during an exam and its effect on test anxiety). In this section, you find out what the signed rank test is and how it’s carried out, and I walk you through an application involving the test of a weight-loss program.

A limitation of the sign test

The sign test has the advantage of being very simple and easy to do by hand. However, because it only looks at whether a value is above or below the median, it doesn’t take the magnitude of the difference into account.

Looking at Tables 17-1 and 17-3, you see that for each data value, the test statistic for the sign test only counts whether or not each data value is greater than or equal to the median in the null hypothesis, m_o . It doesn’t count how great those differences are. For example, in Table 17-3, you can see that the sixth pair had a huge reduction in test anxiety when chewing gum (from 7 down to 2), but the first pair had a very small reduction in test anxiety (from 10 down to 9). Yet both of these differences received the same outcome (a minus sign) in the test statistic for the sign test.

Because it doesn’t take into account how much the values in the data differ from the median, the sign test is less powerful (meaning less able to detect when H_0 is false) than it could be. So if you want to test the median and you want to take the magnitude of the differences into account (and you’re willing to jump through some math hoops to get there), you can conduct the *signed rank test*, also known as the *Wilcoxon signed rank test*. The next section walks you through it.

Stepping through the signed rank test

Just like the sign test, the only condition of the signed rank test is that the data are ordinal or quantitative.

Following are the steps for carrying out the signed rank test on paired data:

1. Set up your hypotheses.

The null hypothesis is $H_0: m = 0$. Your choices for an alternative hypothesis are $H_a: m \neq 0$; $H_a: m > 0$; or $H_a: m < 0$, depending on whether you want to detect any difference, a positive difference, or a negative difference in the pairs.

2. Collect a random sample of paired data.

3. For each observation, calculate the difference for each pair of observations.

4. Calculate the absolute value of each of the differences.

5. Rank the absolute values from smallest to largest.

If two of the absolute values are tied, give each one the average rank of the two values. For example, if the fourth and fifth numbers in order are tied, give each one the rank of 4.5.

6. Add up the ranks that correspond to those original differences from step three that are positive.

The sum of the positive differences is your signed rank test statistic, denoted by SR.

7. Find the *p*-value.

Look at all possible ways that the absolute differences could have appeared in a sample, with either plus or minus signs, assuming that H_0 is true. Find all their test statistics (SR values) from all these possible arrangements by using steps four through six, and compare your SR value to those. The percentage of SR values that are at or beyond your test statistic is your *p*-value.

Minitab can do this step for you.

8. Make your conclusion.

If the *p*-value is less than the predetermined α (typically 0.05), reject H_0 and conclude the median difference is not zero. Otherwise, you can't reject H_0 .



To conduct the Wilcoxon signed rank test using Minitab, enter the differences from step three in a single column. Go to Stat>Nonparametrics>1-Sample Wilcoxon. Click on the name of the variable for your differences in the left-hand box, and it appears in the right-hand Variables box. Click on the circle that says Test Median, and indicate which H_a you want (> 0 , < 0 , or \neq). Click OK, and your test is done. (Note that Minitab calculates the test statistic for the signed rank test a little differently than what you'd get by hand, although the results are close. The reason for the slight calculation difference is beyond the scope of this book.)

How do you handle situations in which a piece of data is exactly equal to the median? Most of the time (including all data sets you will encounter) this occurrence is rare and can be handled by ignoring those data values and reducing the sample size by one for each time the matchup occurs.

Losing weight with signed ranks

This section shows you the signed rank test in action. I first show you each step as if you were doing the process by hand. Then you see the results in Minitab.

Suppose you want to test whether or not a weight-loss plan is effective. You want to look at the median weight loss for people on the plan by using a matched-pairs experiment. You want the magnitude of weight loss to factor into the analysis, which means you use a signed rank test to analyze the data. Here are the steps you follow to conduct the test in this example:

1. Set up your hypotheses. Test $H_0: m = 0$, where m represents the median weight loss (before the program versus after the program). Your alternative hypothesis is $H_a: m > 0$, indicating the median difference in weight loss is positive.
2. Take a random sample of, say, three people and measure them before and after an eight-week weight-loss program. You calculate the difference in weight of each person (weight before the program minus weight after the program). A positive difference means the person lost weight, and a negative difference means they gained weight.

Table 17-4 shows the data and relevant statistics for the weight-loss signed rank test. (Note that I have only three people in this study; this is for illustrative purposes only.) You can see the differences in weight (before – after) in column four.

Table 17-4 Data on Weight Loss Before and After Program

Person	Before	After	Difference	Difference	Rank
1	200	205	-5	5	1
2	180	160	+20	20	2*
3	134	110	+24	24	3*

* Represents ranks associated with a positive difference in weight loss

3. Take the absolute values of the differences. You can see those in column five of Table 17-4.
4. Rank the absolute differences. Column six reflects the ranks of those absolute values, from 1 to 3.
5. Find your test statistic, which is the sum of the ranks corresponding to positive differences. (In other words, you only count ranks of people who lost weight.) For this data set, those ranks you can count are indicated by * in Table 17-4. The sum turns out to be $2 + 3 = 5$. This number, 5, is your test statistic; you can call it SR to designate the signed rank test statistic.
6. Calculate the p -value. Now you need to compare that test statistic to some distribution to see where it stands. To do this, you determine all the possible ways that the three absolute differences (column five of Table 17-4) — 5, 20, and 24 — could have appeared in a sample, with their actual differences taking on plus signs or minus signs. (Assume H_0 is true and the actual differences have a 50-percent chance of being positive or negative, like the flip of a coin.)

Then you find all their test statistics (SR values) from all these possible

arrangements, and compare your SR value, 5, to those. The percentage of the other SR values that are at or beyond your test statistic is your p -value.

For the weight-loss example, you have eight possible ways that you can have absolute differences of 5, 20, and 24 by including either plus or minus signs on each difference (two possible signs for each equals $2 * 2 * 2 = 8$). Those eight possibilities are listed in separate columns of Table 17-5. SR denotes the sum of the positive ranks in each case (these are the test statistics for each possible arrangement).

Table 17-5 Possible Samples with Absolute Differences of 5, 20, and 24

1	2	3	4	5	6	7	8	Rank of $ Diff $
5	-5*	5	5	-5*	-5*	5	-5*	1
20	20	-20*	20	-20*	20	-20*	-20*	2
24	24	24	-24*	24	-24*	-24*	-24*	3
SR = 6	SR = 5	SR = 4	SR = 3	SR = 3	SR = 2	SR = 1	SR = 0	—

* Denotes negative differences

To make sense of Table 17-5, consider the following: The three absolute differences you have in your data set are 5, 20, and 24, which have ranks 1, 2, and 3, respectively (which you can see in Table 17-4). You can find the eight different combinations of 5, 20, and 24 that exist, where you can put either a minus or plus sign on any of those values. For each scenario, I found the signed rank statistic by summing the ranks for only those differences that are positive (the person lost weight). Those ranks are the column nine values in Table 17-5 for data values without an asterisk (*).

For example, column seven has two negative differences, -20 and -24, and one positive difference of 5 (whose rank among the absolute differences is 1 because it's the smallest; see column nine). Summing the positive ranks in column seven produces a signed rank statistic (SR) of 1 because 5 is the only positive number. (You can see in column two the data that you actually observed in the sample.)

Now compare the test statistic, 5 (from step five), to all the values of SR in the last row of Table 17-5. Because you're using $H_a: m > 0$, you can find the percentage of signed ranks (SR) that are at or above the value of 5. You have two of them out of eight, so your p -value (the percentage of possible test statistics beyond or the same as yours if H_0 were true), is $\frac{2}{8} = 0.25$, or 25 percent.

7. Because the p -value (0.25) is greater than the predetermined value of α (typically 0.05), you can't reject H_0 , and you can't say there's positive weight loss via this program. (**Note:** With a sample size of only three, it's difficult to find any real difference, so the weight-loss program may actually be working and this small data set just couldn't determine that, and one person actually gained weight, which doesn't help.)

Figure 17-3 shows the Minitab output for this test, using the data from Table 17-4. The p -value turns out to be 0.211 due to a slight difference in the way that Minitab calculates the test statistic. Note the estimated median found in Figure 17-3 refers to a calculation

made over all possible samples and the medians you would get from them.

Figure 17-3:
Computer
output for
signed rank
test of
weight-loss
data.



Wilcoxon Signed Rank Test: Wt loss

Test of median = 0.000000 versus median > 0.000000

	N	for	Wilcoxon	Estimated
	N	Test	Statistic	P
Wt loss	3		5.0	0.211
				14.75

You also can use the SR statistic to estimate the median of one population (or the median of the difference in a matched-pairs situation). To find a $1 - \alpha$ percent confidence interval for the median using Minitab based on the signed rank test, enter your data into a single column. (If your data represents differences from a matched-pairs data set, enter those differences as one column.) Go to Stat>Nonparametrics>1-Sample Wilcoxon. Click on the name of the variable in the left-hand column, and it appears in the Variables column on the right-hand side. Click the circle that says Confidence Interval, and type in the value of $1 - \alpha$, your confidence level. Click OK.

Pulling Rank with the Rank Sum Test

In This Chapter

- ▶ Comparing two populations by using medians, not means
 - ▶ Conducting the rank sum test
-

In Stats I, when you want to compare two populations, you conduct a hypothesis test for two population means. The most common tool for comparing population means is a *t*-test (see Chapter 3). However, a *t*-test has the condition that the data come from a normal distribution. When conditions for parametric procedures (ones involving normal distributions) aren't met, a nonparametric alternative is always there to save the day.

In this chapter, you work with a nonparametric test that compares the centers of two populations — the *rank sum test*. This test focuses on the *median*, which is the measure of center that's most appropriate in situations where the data isn't symmetric.

Conducting the Rank Sum Test

This section addresses the conditions for the rank sum test and walks you through the steps for conducting the test. You can put your understanding and skill to the test (pun intended) in the section “Performing a Rank Sum Test: Which Real Estate Agent Sells Homes Faster?” later in this chapter.

Checking the conditions

Before you can think about conducting the rank sum test to compare the medians of two populations, you have to make sure your data sets meet the conditions for the test. The conditions for the rank sum test are the following:

- ✓ **The two random samples, one taken from each population, are independent of each other.**

You take care of the first condition in the way you collect your data. Just make sure you aren't using matched pairs, for example, using data from the same person in a pretest and post-test manner. Then the two sets of data would be dependent.

- ✓ **The two populations have the same distribution — that is, their histograms have the same shape.**

You can check this condition by making histograms to compare the shapes of the sample data from the two populations. (See your Stats I textbook or my book *Statistics For Dummies* (Wiley) for help making histograms.)

- ✓ The two populations have the same variance, meaning that the amount of spread in the values is the same.

You can check this condition by finding the variances or standard deviations of the two samples. They should be close. (A hypothesis test for two variances actually exists, but that's outside the scope of this book.)

Notice that the centers of the two populations need not be equal; that's what the test is going to decide.



More sophisticated methods for checking the second and third conditions listed here fall outside the scope of this book. However, checking the conditions as I describe them allows you to find and stay clear of any major problems.

Stepping through the test

The *rank sum test* is a test for the equality of the two population medians — call them η_1 and η_2 . After you've checked the conditions for using the rank sum test (see the preceding section), you conduct the test by following these steps. (**Note:** Minitab can run this test for you, but you should still know what it's doing behind the scenes.)

1. Set up $H_0: \eta_1 = \eta_2$ versus $H_a: \eta_1 > \eta_2$ (a one-sided test); $H_a: \eta_1 < \eta_2$ (a one-sided test); or $H_a: \eta_1 \neq \eta_2$ (a two-sided test), depending on whether you're looking for a positive difference, a negative difference, or any difference between the two population medians.
2. Think of the data as one combined group and assign overall ranks to the values from lowest (rank = 1) to highest.

In the case of ties, give both values the average of the ranks they normally would have been given. For example, suppose the third and fourth numbers (in order) are the same. If the two numbers were different, they would get ranks of 3 and 4, respectively. But because they're the same, you give them the same rank, 3.5, which is the average of 3 and 4. Note that the next number (in order) is the fifth number, which receives rank 5.

3. Sum the ranks assigned to the sample that has the smallest sample size; call this statistic T .

You use the smallest sample in this step according to convention — statisticians like to be consistent. If the sample sizes are equal, sum the ranks for the first sample to get T . If the value of T is small (relative to the total sum of all the ranks from both data sets), the numbers from the first sample tend to be smaller than the second sample, hence the median of the first population may be smaller than the median of the second one.

4. Look at Table A-4(a) and (b) in the appendix, the rank sum tables. For the chosen table, find the column and row for the sample sizes of group one and two,

respectively.

You see two critical values, T_L (the lower critical value) and T_U (the upper critical value). These critical values are the boundaries between rejecting H_0 and not rejecting H_0 .

5. Compare your test statistic, T , to the critical values in Table A-4 in the appendix to conclude whether you can reject H_0 — that the population medians are different.



The method you use to compare these values depends on the type of test you're conducting:

- **One-sided test (Ha has a $>$ or $<$ sign in it):** Table A-4 in the appendix shows the critical values for α level 0.05. For a right-sided test (that means you have $H_a: \eta_1 > \eta_2$), reject H_0 if $T \geq T_U$. For a left-sided test (that means where $H_a: \eta_1 < \eta_2$), reject H_0 if $T \leq T_L$. If you reject H_0 , conclude that the population medians are different and that one of them is greater than the other depending on H_a . (Otherwise you can't conclude that there's a difference in their medians.)
- **Two-sided test:** Table A-4 in the appendix shows the critical values for α level 0.025. Reject H_0 if T falls outside of the interval (T_L, T_U) ; that is, reject H_0 if $T \leq T_L$ or $T \geq T_U$. Conclude that the population medians are not equal. (Otherwise you can't conclude that there's a difference in their medians.)



To conduct a rank sum test in Minitab, enter your data from the first sample in Column 1 and your data from the second sample in Column 2. Go to Stat>Nonparametrics>Mann-Whitney. Click on the name of your Column 1 variable; it appears in the First Sample box. Click on the name of your Column 2 variable; it appears in the Second Sample box. Under Alternative, there's a pull-down menu to select whether your H_a is not equal, greater than, or less than (as indicated by your particular problem). Click OK, and the test is done.

Stepping up the sample size

After the sample sizes reach a certain point, the table values run out. Table A-4 in the appendix (which shows the critical values for rejecting H_0 in the rank sum test) only shows the critical values for sample sizes between three and ten. If both sample sizes are larger than ten, you use a two-sample Z-test to get an approximation for your answer. That's because for large sample sizes the test statistic T for the rank sum test resembles a normal distribution. (So why not use it? As long as you can leave the proof to the professionals!) The larger the two sample sizes are, the better the approximation will be.

So if both sample sizes are more than ten, you conduct steps one through three of the rank sum test as before. Then, instead of looking up the value of T on Table A-4 in the appendix in step four of the rank sum test, you change it to a Z-value (a value on the

standard normal distribution) by subtracting its mean and dividing by its standard error.

$$Z = \frac{T - \frac{n_1(n_1 + n_2 + 1)}{2}}{\sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}}$$

The formula you use to get this Z -value for the test statistic is $Z = \frac{T - \frac{n_1(n_1 + n_2 + 1)}{2}}{\sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}}$, where T is given by step three in the previous section, n_1 is the sample size for the first data set (taken from the first population), and n_2 is the sample size for the second data set (taken from the second population). After you have the Z -value, follow the same procedures that you do for any test involving a Z -value, such as the test for two population means. That is, find the p -value by looking up the Z -value on a Z -table (see your Stats I text) or look on the bottom row of the t -table (which you can find in the appendix), and finding the area beyond it. (If the test is a two-sided test, double the p -value.) If your p -value is less than α , reject H_0 . Otherwise fail to reject H_0 .



In the case where n is large and you use a Z -value for the test statistic, you can still use Minitab (in fact, that's recommended in order to save you the tedium of working through a big example by hand). Refer to the earlier section "Stepping through the test" for the Minitab directions.

Performing a Rank Sum Test: Which Real Estate Agent Sells Homes Faster?

Suppose you want to choose a real estate agent to sell your house, and two agents are in your area. Your most important criteria is to get the house sold fast, so you decide to find out whether one agent sells homes faster. You choose a random sample of eight homes each agent sold in the last year, and for each home, you record the number of days it was on the market before being sold. You can see the data in Table 18-1.

Table 18-1 Market Time for Homes Sold by Two Real Estate Agents

	<i>Agent Suzy Sellfast</i>	<i>Agent Tommy Nowait</i>
House 1	48 days	109 days
House 2	97 days	145 days
House 3	103 days	160 days
House 4	117 days	165 days
House 5	145 days	185 days
House 6	151 days	250 days
House 7	220 days	251 days
House 8	300 days	350 days

Check out the data summarized in *boxplots* (a graph summarizing the data by showing its minimum, first quartile, median, third quartile, and maximum values) in Figure 18-1a and the descriptive statistics in Figure 18-1b. In the following sections, you use this data to see the rank sum test in action. Prepare to be amazed.



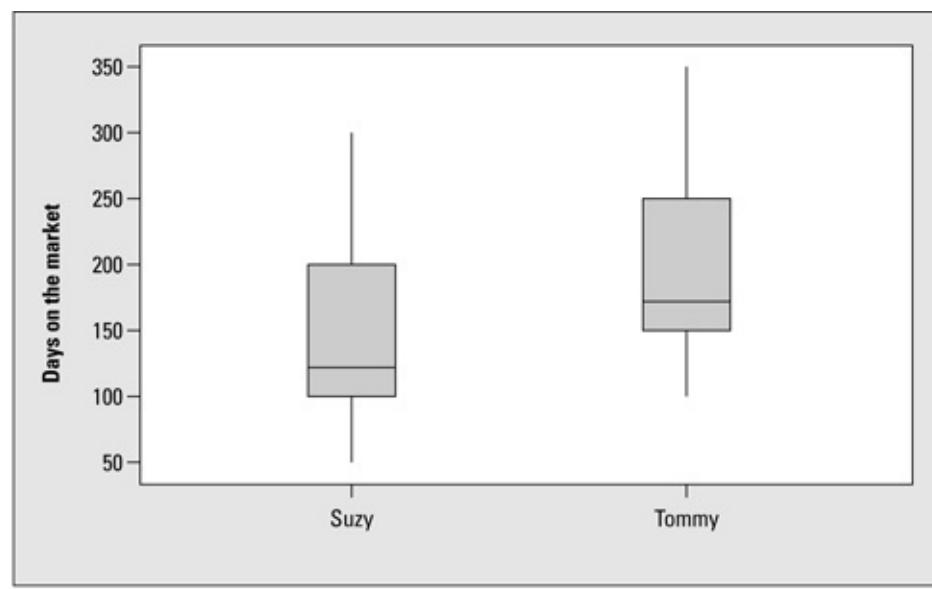
To make two boxplots side by side in Minitab, go to Graph>Boxplots>Simple Multiple Y's. Click on each of your two variables in the left-hand box; they'll appear in the right-hand Variables box. Click OK.

Checking the conditions for this test

In checking the conditions, you know that the data from the two samples are independent, assuming that Suzy and Tommy are competitors.

The boxplots in Figure 18-1a show the same basic shape and amount of variability for each data set. (You don't have enough data to make histograms to check this further.) So based on this data, it isn't unreasonable to assume that the two population distributions of days on the market are the same for the two agents. In Figure 18-1b, the sample standard deviations are close: 79.2 days for Suzy and 77.4 days for Tommy. Because the data meet the conditions for the rank sum test, you can go ahead and apply it to analyze your data.

Figure 18-1:
Boxplots and
descriptive
statistics for
real estate
agent data.



a

Descriptive Statistics: Suzy, Tommy

Variable	Total	Count	Mean	StDev	Minimum	Median	Maximum
Suzy	8	8	147.6	79.2	48.0	131.0	300.0
Tommy	8	8	201.9	77.4	109.0	175.0	350.0

b



To find descriptive statistics (such as the standard deviation) in Minitab, go to Stat>Basic Statistics>Display Descriptive Statistics. Click on Options. Click on the box for each statistic you want to calculate. If a box is checked for a statistic you don't want, click on it again and the check mark disappears.



Figure 18-1b shows that the median for Suzy (131 days on the market) is less than the median for Tommy (175 days). So, on the surface you may say that Suzy sells houses faster, and that's that. But the median doesn't tell the whole story. Looking at Figure 18-1, you see a portion of the two boxplots overlap each other. This means some of the selling times for Suzy and Tommy were close. There's also a great deal of variability in the selling times for each agent, as you can see from the range of values in each boxplot. This tells you that the evidence isn't entirely clear cut. While Suzy may actually be the fastest agent, you can't tell for sure by comparing the boxplots from two samples. You need a hypothesis test to make that final determination.

Testing the hypotheses

The null hypothesis for the real estate agent test is $H_0: \eta_1 = \eta_2$, where η_1 = the median number of days on the market for the population of all Suzy's homes sold in the last year, and η_2 = the median number of days on the market for the population of all Tommy's homes sold in the last year. The alternative hypothesis is $H_a: \eta_1 \neq \eta_2$.



Suppose you looked at the data and developed a hunch that if one of the agents sold homes faster, it was Suzy. However, before you saw the data, you had no preconceived notion as to who was faster. You must base your H_0 and H_a on what your thoughts were *before* you looked at the data, not after. Setting up your hypotheses after you collect the data is unfair and unethical.

After you determine your H_0 and H_a , the time comes to test your data.

Combining and ranking

The first step in the data analysis is to combine all the data and rank the days on the market from lowest (rank = 1) to highest. You can see the overall ranks for the combined data in Table 18-2.

In the case of ties, you give both values the average of the ranks they normally would have received. You can see in Table 18-2 that two values of 145 are in the data set. Because they represent the sixth and seventh numbers in the ordered data set, you give

each of them the same rank of $(6 + 7) \div 2 = 6.5$.

Table 18-2 Ranks of Combined Data from the Real Estate Example

Agent Suzy Sellfast	Overall Rank	Agent Tommy Nowait	Overall Rank
48 days	1	109 days	4
97 days	2	145 days	6.5
103 days	3	160 days	9
117 days	5	165 days	10
145 days	6.5	185 days	11
151 days	8	250 days	13
220 days	12	251 days	14
300 days	15	350 days	16

Finding the test statistic

After you've ranked your data, you determine which group is group one so you can find your test statistic, T . Because the sample sizes are equal, let group one be Suzy because her data is given first. Now sum the ranks from Suzy's data set. The sum of Suzy's ranks is $1 + 2 + 3 + 5 + 6.5 + 8 + 12 + 15 = 52.5$; this value of T is your rank sum test statistic.

Determining whether you can reject H_0

Suppose you want to use an overall α level of 0.05 for this test; using this cutoff means that you use Table A-4(a) in the appendix, because you have a two-sided test at level $\alpha = 0.05$ with 0.025 on each side. Go to the column for $n_1 = 8$ and the row for $n_2 = 8$. You see $T_L = 49$ and $T_U = 87$. You reject H_0 if T is outside this range; in other words, reject H_0 if $T \leq T_L = 49$ or if $T \geq T_U = 87$. Your statistic $T = 52.5$ doesn't fall outside this range; you don't have enough evidence to reject H_0 at the $\alpha = 0.05$ level. So you can't say there's a statistical difference in the median number of days on the market for Suzy and Tommy.

These results may seem very strange given the fact that the medians for the two data sets were so different: 131 days on the market for Suzy compared to 175 days on the market for Tommy. However you have two strikes against you in terms of being able to find a real difference here:

- ✓ **The sample sizes are quite small (only eight in each group).** A small sample size makes it very hard to get enough evidence to reject H_0 .
- ✓ **The standard deviations are both in the high 70s, which is quite large compared to the medians.**

Both of these problems make it hard for the test to actually find anything through all the variability the data show. In other words, it has low power (see Chapter 3).



To conduct the rank sum test by using Minitab, click on

Stat>Nonparametric>Mann-Whitney. Select your two samples and choose your alternate Ha as $>$, $<$, or \neq . The Confidence Level is equal to one minus your value of α . After you make all these settings, click OK.

Figure 18-2 shows the Minitab output when you conduct the rank sum test, or Mann-Whitney test, on the real estate data. To interpret the results in Figure 18-2, you must note that Minitab writes η rather than η for the medians. The results at the bottom of the output say that the test for equal (versus nonequal) medians is significant at the level 0.1149, when adjusting for ties. This is your p -value adjusted for ties. (If no ties are present in your data, you use the results just above that line. That gives you the p -value not adjusted for ties.)

To make your final conclusion, compare your p -value to your predetermined level of α (typically 0.05.) If your p -value is at or below 0.05, you reject H_0 ; otherwise you can't. In this case, because 0.1149 is greater than 0.05, you can't reject H_0 . That means you don't have enough evidence to say the population medians for days on the market for Suzy's versus Tommy's houses are different based on this data. These results confirm your conclusions from the previous section.

Figure 18-2:
Using the
rank sum test
to figure out
who sells
homes faster.

Mann-Whitney Test and CI: Suzy, Tommy		
	N	Median
Suzy	8	131.0
Tommy	8	175.0
 Point estimate for $\eta_{A1}-\eta_{A2}$ is -49.0 95.9 Percent CI for $\eta_{A1}-\eta_{A2}$ is (-137.0, 36.0) $W = 52.5$ Test of $\eta_{A1} = \eta_{A2}$ vs $\eta_{A1} \neq \eta_{A2}$ is significant at 0.1152 The test is significant at 0.1149 (adjusted for ties)		



The Minitab output in Figure 18-2 also provides a confidence interval for the difference in the medians between the two populations, based on the data from these two samples. The difference in the sample medians ($Suzy - Tommy$) is $131.0 - 175.0 = -44.0$. Adding and subtracting the margin of error (these calculations are beyond the scope of this book), Minitab finds the confidence interval for the difference in medians ($Suzy - Tommy$) is $-137.0, +36.0$; the difference in the population medians could be anywhere from -137.0 to 36.0 . Because 0, the value in H_0 , is in this interval, you can't reject H_0 in this case. So again, you can't say that the medians are different, based on this (limited) data set.

Using a rank sum test to compare judges' scoring practices

You can use rank sum tests to compare two groups of judges of a competition to see whether there's a difference in their scores. For example, in figure-skating competitions, the gender of the judges is sometimes suspected to play a role in the scores they give to certain skaters. Suppose you have a men's figure-skating competition with ten judges: five males and five females. You want to know

whether male and female judges score the competitors in the same way, so you do a rank sum test to compare their median scores. Your hypotheses are H_0 : Male and female judges have the same median score versus H_a : They have different median scores. For your sample, you let each judge score every individual. You rank their scores in order from lowest to highest and label M for a male judge and F for a female judge. Your results are the following: F, M, M, M, M, F, F, F, F, M. The value of the test statistic T is the sum of the ranks for group one (the males), which gives you $T = 2 + 3 + 4 + 5 + 10 = 24$. Now compare that to the critical values in Table A-4 in the appendix, where both sample sizes equal five, and you get $T_L = 18$ and $T_U = 37$. Because your test statistic, $T = 24$, is inside this interval, you fail to reject H_0 : Judging is the same for male and female judges. In this situation you don't have enough evidence to say that they differ.

Chapter 19

Do the Kruskal-Wallis and Rank the Sums with the Wilcoxon

In This Chapter

- ▶ Comparing more than two population medians with the Kruskal-Wallis test
 - ▶ Determining which populations are different by using the Wilcoxon rank sum test
-

Statisticians who are in the nonparametrics business make it their job to always find a nonparametric equivalent to a parametric procedure (one that doesn't depend on the normal distribution). And in the case of comparing more than two populations, these stats superheroes don't let us down. In this chapter, you see how the Kruskal-Wallis test works to compare more than two populations as a nonparametric procedure versus its parametric counterpart, ANOVA (see Chapter 9). If Kruskal-Wallis tells you at least two populations differ, this chapter also helps you figure out how to use the Wilcoxon rank sum test to determine which population is different in the same way multiple comparison procedures follow ANOVA (see Chapter 10).

Doing the Kruskal-Wallis Test to Compare More than Two Populations

The Kruskal-Wallis test compares the medians of several (more than two) populations to see whether or not they're different. The basic idea of Kruskal-Wallis is to collect a sample from each population, rank all the combined data from smallest to largest, and then look for a pattern in how those ranks are distributed among the various samples. For example, if one sample gets all the low ranks and another sample gets all the high ranks, perhaps their population medians are different. Or if all the samples have an equal mix of all the ranks, perhaps the medians of the populations are all deemed to be the same. In this section, you see exactly how to conduct the Kruskal-Wallis test using ranks and sums and all that good stuff, and you see it applied to an example comparing airline ratings.

Suppose your boss flies a lot, and she wants you to determine which of three airlines gets the best ratings from customers. You know that ratings involve data that's just not normal (pun intended), so you opt to use the Kruskal-Wallis test. You take three random samples of nine people each from three different airlines. You ask each person to rate his satisfaction with their one airline. Each person uses a scale from 1 (the worst) to 4 (the best). You can see the data from your samples in Table 19-1.



You may be thinking of using ANOVA, the test that compares the means of several populations (see Chapter 9), to analyze this data. But the data from each airline consist of ratings from 1 to 4, which blows the strongest condition of ANOVA — the data from each population must follow a normal distribution. (A *normal distribution* is continuous, meaning it takes on all real numbers in a certain range. Data that are whole numbers like 1, 2, 3, and 4 don't fall under this category.) But don't sweat; a nonparametric alternative fits the bill. The Kruskal-Wallis test compares the medians of several (more than two) populations to see whether they're all the same or not. In other words, it's like ANOVA except that it's done with medians not means.

Table 19-1 Customer Ratings of Three Airlines

Airline A Rating	Airline B Rating	Airline C Rating
4	2	2
3	3	3
4	3	3
4	3	2
3	4	2
3	4	1
2	3	3
3	4	2
4	3	2

In looking at the data in Table 19-1, it appears that airlines A and B have better ratings than airline C. However, the data have a lot of variability, so you have to conduct a hypothesis test before you can make any general conclusions beyond this data set.

In this section, you discover how to check the conditions of the Kruskal-Wallis test, set it up, and carry it out step by step.

Checking the conditions

The following conditions must be met in order to conduct the Kruskal-Wallis test:

- ✓ The random samples taken from each population are independent. (This means matched-pairs data like in Chapter 17 are out of this picture.)
- ✓ All the populations have the same distribution, meaning their shapes are the same as seen on a histogram. (Note they don't specify what that distribution is.)
- ✓ The variances of the populations are the same. The amount of spread in the population values is the same from one population to the next.

Note that these conditions mention shape and spread, but not the center of the distributions. The test is trying to determine whether the populations are centered at the same place.



In nonparametrics, you often see the word *location* used in reference to a population distribution rather than the word *center*, although the two words mean about the same thing. Location indicates where the distribution is sitting on the number line. If you have two bell-shaped curves with the same variance and one has mean 10 and the other has mean 15, the second distribution is located 5 units to the right of the first. In other words, its location is a 5-unit shift to the right of the first distribution. In nonparametrics, where you don't have bell-shaped distributions, you typically use the median as a measure of location of a distribution. So throughout this discussion, you could use the word *median* instead of location (although location leaves it a bit more open).

Regarding the airline survey, you know that the samples are independent, because you didn't use the same person to rate more than one airline. The other two conditions have to do with the distributions the samples came from; each population must have the same shape and the same spread. You can examine both conditions by looking at boxplots of the data (see Figure 19-1) and descriptive statistics, such as the median, standard deviation, and the rest of the summary statistics making up the boxplots (see Figure 19-2).

The boxplots in Figure 19-1 all have the same shape, and their standard deviations, shown in Figure 19-2, are very close. All this evidence taken together allows you to go ahead with the Kruskal-Wallis test. (Looking at the overlap in the boxplots for airlines A and B in Figure 19-1, you also can make an early prediction that airlines A and B have similar ratings. Whether C is different enough from A and B is impossible to say without running the hypothesis test.)

Figure 19-1:
Boxplots
comparing
the ratings of
three airlines.

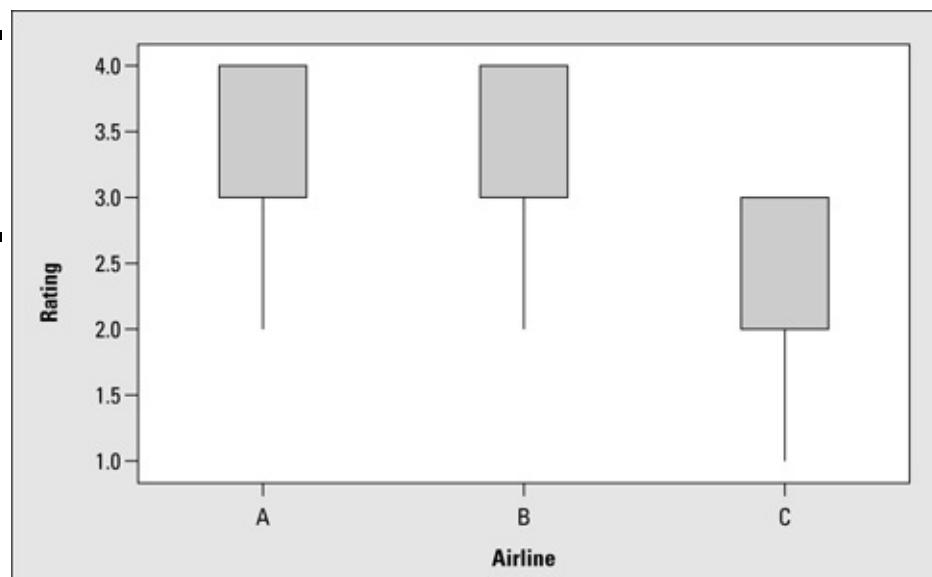


Figure 19-2:
Descriptive statistics comparing the ratings of three airlines.

Descriptive Statistics: Rating

Variable	Airline	StDev	Minimum	Q1	Median	Q3	Maximum
Rating	A	0.707	2.000	3.000	3.000	4.000	4.000
	B	0.667	2.000	3.000	3.000	4.000	4.000
	C	0.667	1.000	2.000	2.000	3.000	3.000



Either a boxplot or a histogram can tell you about the shape and spread of a distribution (as well as the center). The *boxplot* is a common type of graph to use for nonparametric procedures because it displays the median (the nonparametric statistic of choice) rather than the mean. At its best, a *histogram* shows the shape of the data; it doesn't directly tell where the center is — you just have to eyeball it.



Note that the boxplots in Figure 19-1 don't have lines through them, as you may expect. That's because in each case the median happens to equal Q1, the first quartile (see Figure 19-2). This situation can happen with rank data when many ranks take on the same value.



To make boxplots of each sample of data show up side by side on one graph (called *side-by-side boxplots*, cleverly) in Minitab, click on Graph>Box Plots and select the Multiple Y's Simple version. In the left-hand box, click on each of the column names for your data sets. They each appear in the Graph Variables window on the right. Click OK and you get a set of boxplots that are side by side, all on the same graph using the same scale (slick, huh?).

Setting up the test

The Kruskal-Wallis test assesses H_0 : All k populations have the same location versus H_a : The locations of at least two of the k populations are different. Here, k is the number of populations you're comparing.



In H_0 , you see that all the populations have the same location (which means they all sit on top of each other on the number line and are in essence the same population). H_a is looking for the opposite situation in this case. However, the opposite of "the locations are all equal" isn't "the locations are all different." The opposite is that at least two of them are different. Failure to recognize this difference will lead you to believe all the populations differ when, in reality, there may only be two that differ, and the rest are all the same. That's why you see H_a stated the way it

is in the Kruskal-Wallis test. (The same idea holds for comparing means using ANOVA; see Chapter 9.)

For the airline satisfaction example (see Table 19-1), your setup looks like this: Ho: The satisfaction ratings of all three airlines have the same median versus Ha: The median satisfaction ratings of at least two airlines are different.

Conducting the test step by step

After you've determined your hypotheses and checked the conditions, you can carry out the test. Here are the steps for conducting the Kruskal-Wallis test using the airline example to show how each step works:

1. Rank all the numbers in the entire data set from smallest to largest (using all samples combined); in the case of ties, use the average of the ranks that the values would have normally been given.

Figure 19-3 shows the results for ranking and summing the data in the airline example; you can see how to rank the ties. For example, you have only one 1, which is given rank 1. Then you have seven 2s, which normally would have gotten ranks 2, 3, 4, 5, 6, 7, and 8. Because the 2s are all equal, you give each of them the average of all these

ranks, which is $\frac{(2+3+4+5+6+7+8)}{7} = 5$. Similarly, you see twelve 3s, whose ranks would be 9 through 20. Because they're all equal, give them each a rank equal to $\frac{(9+10+\dots+20)}{12} = 14.5$. Finally, you see seven 4s, each with rank 24, which is the average of their would-be ranks, ranging from 21 to 27.

Figure 19-3:
Rankings and
rank sum for
the airline
example.

Airline A		Airline B		Airline C	
Rating	Rank	Rating	Rank	Rating	Rank
4	24	2	5	2	5
3	14.5	3	14.5	3	14.5
4	24	3	14.5	3	14.5
4	24	3	14.5	2	5
3	14.5	4	24	2	5
3	14.5	4	24	1	1
2	5	3	14.5	3	14.5
3	14.5	4	24	2	5
4	24	3	14.5	2	5
$T_1 = 159$		$T_2 = 149.5$		$T_3 = 69.5$	

2. Total the ranks for each of the samples; call those totals T_1 , T_2 , . . . , T_k , where k is the number of populations.

The totals of the ranks in each column of Figure 19-3 are $T_1 = 159$ (the total ranks for airline A), $T_2 = 149.5$, and $T_3 = 69.5$. In the steps that follow, you use these rank totals in the Kruskal-Wallis test statistic (denoted KW). (Note T_1 and T_2 are close to equal, but T_3 is much lower, giving the impression that airline C may be the odd man out.)

3. Calculate the Kruskal-Wallis test statistic, $KW = \frac{12}{n(n+1)} \sum \frac{T_j^2}{n_j} - 3(n+1)$, where n is the total number of observations (all sample sizes combined).

For the airline example, the Kruskal-Wallis test statistic is

$$KW = \frac{12}{27(27+1)} \left(\frac{159^2}{9} + \frac{149.5^2}{9} + \frac{69.5^2}{9} \right) - 3(27+1), \text{ which equals } 0.0159 * 5,829.056 - 3(28) = 8.52.$$

4. Find the p -value for your KW test statistic by comparing it to a Chi-square distribution with $k - 1$ degrees of freedom (see Table A-3 in the appendix).

For the airline example, you look at the Chi-square table (Table A-3 in the appendix) and find the row with $3 - 1 = 2$ degrees of freedom. Then look at where your test statistic (8.52) falls in that row. Because 8.52 lies between 7.38 and 9.21 (shown on the table in row two), the p -value for 8.52 lies between 0.025 and 0.010 (shown in their respective column headings.)

5. Make your conclusion about whether you can reject H_0 by examining the p -value.

You can reject H_0 : All populations have the same location, in favor of H_a : At least two populations have differing locations, if the p -value associated with KW is $< \alpha$, where α is 0.05 (or your prespecified α level). Otherwise, you fail to reject H_0 .

Following the airline example, because the p -value is between 0.010 and 0.025, which are both less than $\alpha = 0.05$, you can reject H_0 . You conclude that the ratings of at least two of the three airlines are different.



To conduct the Kruskal-Wallis test by using Minitab, enter your data in two columns, the first column representing the actual data values and the second column representing which population the data came from (for example, 1, 2, 3). Then click on Stat>Nonparametrics>Kruskal-Wallis. In the left-hand box, click on column one; it appears on the right side as your *response variable*. Then click on column two in the left-hand box. This column appears on the right side as the *factor variable*. Click OK, and the KW test is done. The main results of the KW test are shown in the last two lines of the Minitab output.

The results of the Minitab data analysis of the airline data are shown in Figure 19-4. On the second-to-last line of Figure 19-4, you can see the KW test statistic for the airline example is 8.52, which matches the one you found by hand (whew!). The exact p -value from Minitab is 0.014.

Figure 19-4:
Comparing
ratings of
three airlines
by using the
Kruskal-
Wallis test.

Kruskal-Wallis Test: Rating versus Airline

Kruskal-Wallis Test on Rating

Airline	N	Median	Ave Rank	Z
A	9	3.000	17.7	1.70
B	9	3.000	16.6	1.21
C	9	2.000	7.7	-2.91
Overall	27		14.0	

H = 8.52 DF = 2 P = 0.014
H = 9.70 DF = 2 P = 0.008 (adjusted for ties)

However, this data set has quite a few ties, and the formulas have to adjust a bit for that (in ways that go outside the scope of this book). Taking those ties into account, the computer gives you $KW = 9.70$ with a p -value of 0.008. The total evidence here says the same result loud and clear — reject H_0 : The ratings for the three airlines have the same location. You conclude that the ratings of at least two of the airlines are different. (But which ones? The answer comes in the next section.)

Pinpointing the Differences: The Wilcoxon Rank Sum Test

Suppose you reject H_0 in the Kruskal-Wallis test, meaning you have enough evidence to conclude that at least two of the populations have different medians. But you don't know which ones are different. When you find that a set of populations don't all share the same median, the next question is very likely to be, "Well then, which ones are different?" To find out which populations are different after the Kruskal-Wallis test has rejected H_0 , you can use the *Wilcoxon rank sum test* (also known as the *Mann-Whitney test*).



You can't go looking for differences in specific pairs of populations until you've first established that at least two populations differ (that is, H_0 is rejected in the Kruskal-Wallis test). If you don't make this check first, you can encounter a ton of problems, not the least of which being a much-increased chance of making the wrong decision.

In the following sections, you see how to conduct pairwise comparisons and interpret them in order to find out where the differences lie among the k population medians you're studying.

Pairing off with pairwise comparisons

The rank sum test is a nonparametric test that compares two population locations (for example, their medians). When you have more than two populations, you conduct the

rank sum test on every pair of populations in order to see whether differences exist. This procedure is called conducting *pairwise comparisons* or *multiple comparisons*. (See Chapter 10 for info on the parametric version of multiple comparisons.) For example, because you’re comparing three airlines in the airline satisfaction example (see Table 19-1), you have to run the rank sum test three times to compare airlines A and B, A and C, and B and C. So you need three pairwise comparisons to figure out which populations are different.



To determine how many pairs of comparisons you need if you’re given k

$$\frac{k(k-1)}{2}$$

populations, you use the formula $\frac{k(k-1)}{2}$. You have k populations to choose from first and then $k - 1$ populations left to compare them with. Finally, you don’t care what the order is among the populations (as long as you keep track of them); so you divide by two because you have two ways to order any pair (for example, comparing A and B gives you the same results as comparing B and A). In the airlines example,

you have $k = 3$ populations, so you should have $\frac{k(k-1)}{2} = \frac{3(3-1)}{2} = 3$ pairs of populations to compare, which matches what was determined previously. (For more information and examples on how to count the number of ways to choose or order a group of items by using permutations and combinations, see another book I authored, *Probability For Dummies* [Wiley].)

Carrying out comparison tests to see who’s different

The Wilcoxon rank sum test assesses H_0 : The two populations have the same location versus H_a : The two populations have different locations. Here are the general steps for using the Wilcoxon rank sum test for making comparisons:

1. Check the conditions for the test by using descriptive statistics and histograms for the last two and proper sampling procedures for the first one:

- The two samples must be from independent populations.
- The populations must have the same distribution (shape).
- The populations must have the same variance.

2. Set up your H_0 : The two medians are equal versus H_a : The two medians aren’t equal.

3. Combine all the data and rank the values from smallest to largest.

4. Add up all the ranks from the first sample (or the smallest sample if the sample sizes are not equal).

This result is your test statistic, T .

5. Compare T to the critical values in Table A-4 in the appendix, in the row and column corresponding to the two sample sizes (denoted T_L and T_U).

If T is at or beyond the critical values (less than or equal to the lower one [T_L] or

greater than or equal to the upper one [T_U]), reject H_0 and conclude the two population medians are different. Otherwise, you can't reject H_0 .

6. Repeat steps one through five on every pair of samples in the data set and draw conclusions.

Sort through all the results to see the overall picture of which pairs of populations have the same median and which ones don't.



To conduct the Wilcoxon rank sum test for pairwise comparisons in Minitab, refer to Chapter 18. Note that Minitab calls this test by its other name, the Mann-Whitney test.

You can see the Minitab results of the three Wilcoxon rank sum tests comparing airlines A and B, A and C, and B and C in Figures 19-5a, 19-5b, and 19-5c, respectively.

Figure 19-5a compares the ratings of airlines A and B. The p -value (adjusted for ties) is 0.7325, which is much higher than the 0.05 you need to reject H_0 . So you can't conclude that airlines A and B have satisfaction ratings with different medians. Figure 19-5b shows that the p -value for comparing airlines A and C is 0.0078. Because this p -value is a lot smaller than the typical α level of 0.05, it's very convincing evidence that airlines A and C don't have the same median ratings. Figure 19-5c also has a small p -value (0.0107), which gives evidence that airlines B and C have significantly different ratings.

Figure 19-5:
Wilcoxon
rank sum
tests
comparing
ratings of two
airlines at a
time.

Mann-Whitney Test and CI: Airline A, Airline B

	N	Median
A	9	3.000
B	9	3.000

Point estimate for ETA1-ETA2 is -0.000
95.8 Percent CI for ETA1-ETA2 is (-1.000,1.000)
W = 89.5
Test of ETA1 = ETA2 vs ETA1 not = ETA2 is significant at 0.7573
The test is significant at 0.7325 (adjusted for ties)

a

Mann-Whitney Test and CI: Airline A, Airline C

	N	Median
A	9	3.000
C	9	2.000

Point estimate for ETA1-ETA2 is 1.000
95.8 Percent CI for ETA1-ETA2 is (0.000,2.000)
W = 114.5
Test of ETA1 = ETA2 vs ETA1 not = ETA2 is significant at 0.0118
The test is significant at 0.0078 (adjusted for ties)

b

Mann-Whitney Test and CI: Airline B, Airline C

	N	Median
B	9	3.000
C	9	2.000

Point estimate for ETA1-ETA2 is 1.000
95.8 Percent CI for ETA1-ETA2 is (0.000,2.000)
W = 113.0
Test of ETA1 = ETA2 vs ETA1 not = ETA2 is significant at 0.0171
The test is significant at 0.0107 (adjusted for ties)

c

Examining the medians to see how they're different

Rejecting H_0 for a multiple comparison means you conclude the two populations you examined have different medians. There are two ways to proceed from here to see how the medians differ for all pairwise comparisons:

- ✓ You can look at side-by-side boxplots of all the samples and compare their medians (located at the line in the middle of each box).
- ✓ You can calculate the median of each sample and see which ones are higher and which ones are lower from the populations you have concluded are statistically different.

From the previous section, you see that the pairwise comparisons for the airline data conducted by Wilcoxon rank sum tests conclude that the ratings of airlines A and B aren't found to be different, but both of them are found to be different from airline C.

But you can say even more; you can say how the differing airline compares to the others.

Going back to Figure 19-2, you see the medians of both airlines A and B are 3.0, while the median of airline C is only 2.0. That difference means airlines A and B have similar ratings, but airline C has lower ratings than A and B. The boxplots in Figure 19-1 confirm these results.

Chapter 20

Pointing Out Correlations with Spearman's Rank

In This Chapter

- ▶ Understanding correlation from a nonparametric point of view
- ▶ Finding and interpreting Spearman's rank correlation

Data analysts commonly look for and try to quantify relationships between two variables, x and y . Depending on the type of data you're dealing with in x and y , you use different procedures for quantifying their relationship.

When x and y variables are *quantitative* (that is, their possible outcomes are measurements or counts), the correlation coefficient (also known as the *Pearson's correlation coefficient*) measures the strength and direction of their linear relationship. (See Chapter 4 for all the info on Pearson's correlation coefficient, denoted by r .) If x and y are both *categorical* variables (meaning their possible outcomes are categories that have no numerical meaning, such as male and female), you use Chi-square procedures and conditional probabilities to look for and describe their relationship. (I lay out all that machinery in Chapters 13 and 14.)

Then there's a third type of variable, called *ordinal* variables; their values fall into categories, but the possible values can be placed into an order and given a numerical value that has some meaning, for example, grades on a scale of A = 4, B = 3, C = 2, D = 1, and E = 0 or a student's evaluation of a teacher on a scale from best = 5 to worst = 1. To look for a relationship between two ordinal variables like these, statisticians use *Spearman's rank correlation*, the nonparametric counterpart to Pearson's correlation coefficient (covered in Chapter 4). In this chapter, you see why ordinal variables don't meet Pearson's conditions, and you find out how to use and interpret Spearman's rank correlation to correctly quantify and interpret relationships involving ordinal variables.

Pickin' On Pearson and His Precious Conditions

Pearson's correlation coefficient is the most common correlation measure out there, and many data analysts think it's the *only* one out there. Trouble is, Pearson's correlation has certain conditions that must be met before using it. If those conditions aren't met, Spearman's correlation is waiting in the wings.



The Pearson correlation coefficient r (the correlation) is a number that measures

the direction and strength of the linear relationships between two variables x and y . (For more info on the correlation, see Chapter 4.)

Several conditions have to be met for ol' Pearson:

- ✓ **Both variables x and y must be numerical (or quantitative).** They must represent measurements with no restriction on their level of precision. For example, numbers with many places after the decimal point (such as 12.322 or 0.219) must be possible.
- ✓ **The variables x and y must have a linear relationship (as shown on a scatterplot; see Chapter 4).**
- ✓ **The y values must have a normal distribution for each x , with the same variance at each x .**



Typically with ordinal variables, you won't see many different categories offered or compared for reasons of simplicity. This means there won't be enough numerical values to try to build a linear regression model involving ordinal variables like you can with two quantitative variables.

In another scenario, if you have a gender variable with categories male and female, you can assign the numbers 1 and 2 to each gender, but those numbers have no numerical meaning. Gender isn't an ordinal variable; rather it's a *categorical variable* (a variable that places individuals into categories only). Categorical variables also don't lend themselves to linear relationships, so they don't meet Pearson's conditions either. (To explore relationships between categorical variables, see Chapter 14.)

Who are these guys? A look at the people behind the statistics

Some people are lucky enough to have a statistic actually named after them. Typically, the person who came up with the statistic in the first place, recognizing a need for it and coming up with a solution, gets the honor. If the new statistic gets picked up and used by others, it eventually takes on the name of its inventor.

Spearman's rank correlation is named after its inventor, Charles Edward Spearman (1863–1945). He was an English psychologist who studied experimental psychology, worked in the area of human intelligence, and was a professor for many years at the University College London. Spearman followed closely the work of Francis Galton, who originally developed the concept of correlation. Spearman developed his rank correlation in 1904.

Pearson's correlation coefficient was developed several years prior, in 1893 by Karl Pearson, one of Spearman's colleagues at University College London and another follower of Galton. Pearson and Spearman didn't get along; Pearson had an especially strong and volatile personality and had problems getting along with quite a few people, in fact.

Scoring with Spearman's Rank Correlation

Spearman's rank correlation doesn't require the relationship between the variables x and y to be linear, nor does it require the variables to be numerical. Rather than examining a linear relationship between x and y , Spearman's rank correlation tests whether two ordinal and/or quantitative variables are independent (in other words, not related to each other). **Note:** Spearman's rank applies to ordinal data only. To test to see if two categorical (and nonordinal) variables are independent, you use a Chi-square test; see Chapter 14.



Spearman's rank correlation is the same as Pearson's correlation except that it's calculated based on the *ranks* of the x and y variables (that is, where they stand in the ordering; see Chapter 16) rather than their actual values. You interpret the value of Spearman's rank correlation, r_s , the same way you interpret Pearson's correlation, r (see Chapter 4). The values of r_s can go between -1 and $+1$. The higher the magnitude of r_s (in the positive or negative directions), the stronger the relationship between x and y . If r_s is zero, x and y are independent. And as with r , if the correlation between x and y is not zero, you can't say whether or not they're independent.

In this section, you see how to calculate and interpret Spearman's rank correlation and apply it to an example.

Figuring Spearman's rank correlation

The notation for Spearman's rank correlation is r_s , where s stands for Spearman. To find r_s , you do the steps listed in this section. Minitab does the work for you in steps two through six, although some professors may ask you to do the work by hand (not me of course).

- 1. Collect the data in the form of pairs of values x and y .**
- 2. Rank the data from the x variable where $1 = \text{lowest}$ to $n = \text{highest}$, where n is the number of pairs of data in the data set.**

This step gives you a new set of data for the x variable called the *ranks* of the x values. If any of the values appear more than once, Minitab assigns each tied value the average of the ranks they would normally be given if they weren't tied.

- 3. Complete step two with the data from the y variable.**

This step gives you a new data set called the *ranks* of the y -values.

- 4. Find the standard deviation of the ranks of the x -values, using the usual formula**

for standard deviation, $\sqrt{\frac{\sum(x - \bar{x})^2}{n-1}}$; call it s_{xx} . In a similar manner find the standard

deviation of the ranks of the y -values using $\sqrt{\frac{\sum(y - \bar{y})^2}{n-1}}$; call it s_{yy} .

Note that n is the sample size, \bar{x} is the mean of the ranks of the x values, and \bar{y} is the mean of the ranks of the y values.

5. Find the covariance of the x - y values, using the formula $\text{Cov}(x, y) = \frac{\sum \sum (x - \bar{x})(y - \bar{y})}{n-1}$; call it s_{xy} .

The covariance of x and y is a measure of the total deviation of the x and y values from the point (\bar{x}, \bar{y}) .

6. Calculate the value of Spearman's rank correlation by using the formula $r_s = \frac{s_{xy}}{s_{xx}s_{yy}}$.



The formula for Spearman's rank correlation is just the same as the formula for Pearson's correlation coefficient, except the data Spearman uses is the ranks of x and y rather than the original x - and y -values as used by Pearson. So Spearman just cares about the order of the values of the x 's and the y 's, not their actual values.



To calculate Spearman's rank correlation straightaway by using Minitab, rank the x -values, rank the y -values, and then find the correlation of the ranks. That is, go to Data>Rank and click on the x variable to get x ranks. Then do the same thing to get the y ranks. Go to Stat>Basic Statistics>Correlation, click on the two columns representing ranks, and click OK.

Watching Spearman at work: Relating aptitude to performance

Knowing the process of how to calculate Spearman's rank correlation is one thing, but if you can apply it to real-world situations, you'll be the golden child of the statistics world (or at least your statistics class). So, try putting yourself in this section's scenario to get the full effect of Spearman's rank correlation.

You're a statistics professor, and you give exams every now and then (it's a dirty job, but someone's got to do it). After looking at students' final grades over the years (yes, you're an old professor, or at least in your mid-40s), you notice that students who do well in your class tend to have a better aptitude (background ability) for math and statistics. You want to check out this theory, so you give students a math and statistics aptitude test on the first day of the course; you want to compare students' aptitude test scores with their final grades at the end of the course.

Now for the specifics. Your variables are x = aptitude test score (using a 100-point pretest

on the first day of the course) and y = final grade on a scale from 1 to 5 where 1 = F (failed the course), 2 = D (passed), 3 = C (average), 4 = B (above average), and 5 = A (excellent). The y variable, final grade, is an ordinal variable, and the x variable, aptitude, is a numerical variable. You want to find out whether there's a relationship between x and y . You collect data on a random sample of 20 students; the data are shown in Table 20-1. This is step one of the process of calculating Spearman's rank correlation (from the steps listed in the previous section).

Table 20-1 Aptitude Test Scores and Final Grades in Statistics

<i>Student</i>	<i>Aptitude</i>	<i>Final Grade</i>
1	59	3
2	47	2
3	58	4
4	66	3
5	77	2
6	57	4
7	62	3
8	68	3
9	69	5
10	36	1
11	48	3
12	65	3
13	51	2
14	61	3
15	40	3
16	67	4
17	60	2
18	56	3
19	76	3
20	71	5

Using Minitab, you get a Spearman's rank correlation of 0.379. The following discussion walks you through steps two through six as you do this correlation yourself. This is likely what you may be asked to do on an exam.

Steps two and three of finding Spearman's rank correlation are to rank the aptitude test scores (x) from lowest (1) to highest; then rank the final grades (y) from lowest (1) to highest. Note that the final exam grades have several ties, so you use average ranks. For example, in column three of Table 20-1 you see a single 1 for Student 10, which gets rank 1. Then you see four 2s for Students 2, 5, 13, and 17. Those ranks, had they not been tied,

would have been 2, 3, 4, and 5. The average of these four ranks is $\frac{2+3+4+5}{4} = \frac{14}{4} = 3.5$.

Each of the 2s in column three, therefore, receives rank 3.5.

Table 20-2 shows the original data, the ranks of the aptitude scores (x), and the ranks of the final grades (y) as calculated by Minitab.

Table 20-2 Aptitude Test Scores, Final Exam Grades, and Rank

Student	Aptitude	Rank of Aptitude	Final Grade	Rank of Final Grade
1	59	9	3	10.5
2	47	3	2	3.5
3	58	8	4	17.0
4	66	14	3	10.5
5	77	20	2	3.5
6	57	7	4	17.0
7	62	12	3	10.5
8	68	16	3	10.5
9	69	17	5	19.5
10	36	1	1	1.0
11	48	4	3	10.5
12	65	13	3	10.5
13	51	5	2	3.5
14	61	11	3	10.5
15	40	2	3	10.5
16	67	15	4	17.0
17	60	10	2	3.5
18	56	6	3	10.5
19	76	19	3	10.5
20	71	18	5	19.5

For step four of finding Spearman's rank correlation, you have Minitab calculate the standard deviation of the aptitude test score ranks (located in column two of Table 20-2) and the standard deviation of the final grades (located in column four of Table 20-2). In step five, you have Minitab calculate the covariance of the ranks of aptitude test scores and final grade ranks. These statistics are shown in Figure 20-1.

Figure 20-1:
Standard deviations and covariance of ranks of aptitude (x) and final grade (y).

Descriptive Statistics: Ranks of X, Ranks of Y		
Variable		StDev
Ranks of X		5.92
Ranks of Y		5.50
Covariances: Ranks of X, Ranks of Y		
	Ranks of X	Ranks of Y
Ranks of X	35.0000	
Ranks of Y	12.3421	30.2632

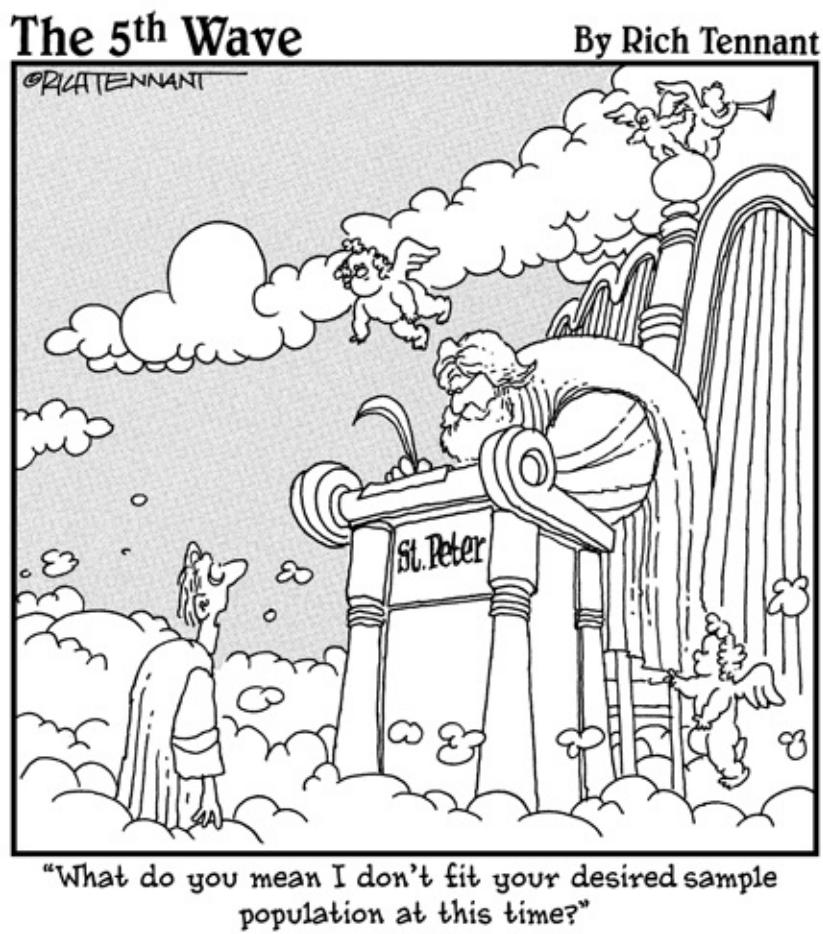
For the sixth and final step of finding Spearman's rank correlation, calculate r_s by taking

the covariance of the ranks of x and y (found in the lower left corner) divided by the standard deviation of the ranks of x (called s_{xx}) times the standard deviation of the ranks of y (called s_{yy}). You get $\frac{12.34}{5.92 \times 5.50} = 0.379$. This matches the value for Spearman's correlation that Minitab found straightaway.

This Spearman's rank correlation of 0.379 is fairly low, indicating a weak relationship between aptitude scores before the course and final grades at the end of the course. The moral of the story? If a student isn't the sharpest tack in the bunch, he can still hope, and if he comes in on top, he may not go out the same way. Although, there's still something to be said about working hard during the course (and buying *Statistics II For Dummies* certainly doesn't hurt!).

Part VI

The Part of Tens



In this part . . .

This part is a staple in *For Dummies* books, and for good reason. In this part, you get useful insight about what statistical conclusions to be wary of, along with the many ways that statistics is used in the workplace. I also tell you why you shouldn't try to escape statistics, but rather plunge right in, because knowing statistics means job security in practically anything you do.

Chapter 21

Ten Common Errors in Statistical Conclusions

In This Chapter

- ▶ Recognizing and avoiding mistakes when interpreting statistical results
- ▶ Deciding whether or not someone's conclusions are credible

Stats II is all about building models and doing data analysis. It focuses on looking at data and figuring out the story behind it. It's about making sure that the story is told correctly, fairly, and comprehensively. In this chapter, I discuss some of the most common errors I've seen as a teacher and statistical consultant for many moons. You can use this "not to do" list to pull ideas together for homework and reports or as a quick review before a quiz or exam. Trust me — your professor will love you for it!

Claiming These Statistics Prove . . .

Be skeptical of anyone who uses *these statistics* and *prove* in the same sentence. The word *prove* is a definitive, end-all-be-all, case-closed, lead-pipe-lock sort of concept, and statistics by nature isn't definitive. Instead, statistics gives you evidence for or against someone's theory, model, or claim based on the data you collected; then it leaves you to your own conclusions. Because the evidence is based on data that changes from sample to sample, the results can change as well — that's the challenge, the beauty, and sometimes the frustration of statistics. The best you can say is that your statistics suggest, lead you to believe, or give you sufficient evidence to conclude — but never go as far as to say that your statistics prove anything.

It's Not Technically Statistically Significant, But . . .



After you set up your model and test it with your data, you have to stand by the conclusions no matter how much you believe they're wrong. Statistics must lend objectivity to every process.

Suppose Barb, a researcher, has just collected and analyzed the heck out of her data, and she still can't find anything. However, she knows in her heart that her theory holds true, even if her data can't confirm it. Barb's theory is that dogs have ESP — in other words, a "sixth sense." She bases this theory on the fact that her dog seems to know when she's leaving the house, when he's going to the vet, and when a bath is imminent because he

gets sad and finds a corner to hide in.

Barb tests her ESP theory by studying ten dogs, placing a piece of dog food under one of two bowls and asking each dog to find the food by pushing on a bowl. (Assume the bowl is thick enough that the dogs can't cheat by smelling the food.) She repeats this process ten times with each dog and records the number of correct responses. If the dogs don't have ESP, you would expect that they would be right 50 percent of the time because each dog has two bowls to choose from and each bowl has an equal chance of being selected.

As it turns out, the dogs in Barb's study were right 55 percent of the time. Now, this percentage is technically higher than the long-term expected value of 50 percent, but it's not enough (especially with so few dogs and so few trials) to warrant statistical significance. In other words, Barb doesn't have enough evidence for the ESP theory. But when Barb presents her results at the next conference she attends, she puts a spin on her results by saying, "The dogs were correct 55 percent of the time, which is more than 50 percent. These results are *technically* not enough to be statistically significant, but I believe they do show some evidence that dogs have ESP" (causing every statistician in the room to scream "NOT!").

Some researchers use this kind of conclusion all the time — skating around the statistics when they don't go their way. This game is very dangerous because the next time someone tries to replicate Barb's results (and believe me, someone always does), they find out what you knew from the beginning (through ESP?): When Barb starts packing to leave the house, her dog senses trouble coming and hides. That's all.

Concluding That x Causes y

Do you see the word that makes statisticians nervous? The first two seem pretty tame, and x and y are just letters of the alphabet, it's got to be that word *cause*. Of all the words used too loosely in statistics, *cause* tops the list.

Here's an example of what I mean. For your final report in stats class, you study which factors are related to a student's final exam score. You collect data on 500 statistics students, asking each one a variety of questions, such as "What was your grade on the midterm?"; "How much sleep did you get the night before the final?"; and "What's your GPA?" You conduct a multiple linear regression analysis (using techniques from Chapter 5) and conclude that study time and the amount of sleep the night before the test are the most-important factors in determining exam scores. You write up all your analyses in a paper, and at the very end you say, "These results demonstrate that more study time and a good night of sleep the night before cause a student's exam grade to increase."

I was with you until you said the word *cause*. You can't say that more sleep or more study time causes an increase in exam score. The data you collected shows that people who get a lot of sleep and study a lot do get good grades, and those who don't do those things don't get the good grades. But that result doesn't mean a flunkie can just sleep and

study more and all will be okay. This theory is like saying that because an increase in height is related to an increase in weight, you can get taller by gaining weight.

The problem is that you didn't take the same person, change his sleep time and study habits, and see what happened in terms of his exam performance (using two different exams of the same difficulty). That study requires a *designed experiment*. When you conduct a *survey*, you have no way of controlling other related factors going on, which can muddy the waters, like quality of studying, class attendance, grades on homework, and so on.



The only way to control for other factors is to do a randomized experiment (complete with a treatment group, a control group, and controls for other factors that may ordinarily affect the outcome). Claiming causation without conducting a randomized experiment is a very common error some researchers make when they draw conclusions.

Assuming the Data Was Normal

The operative word here is *assuming*. To break it down simply, an assumption is something you believe without checking. Assumptions can lead to wrong analyses and incorrect results — all without the person doing the assuming even knowing it.

Many analyses have certain requirements. For example, data should come from a normal distribution (the classic distribution that has a bell shape to it). If someone says, “I assumed the data was normal,” she just assumed that the data came from a normal distribution. But is having a normal distribution an assumption you just make and then move on, or is more work involved? You guessed it — more work.

For example, in order to conduct a one-sample *t*-test (see Chapter 3), your data must come from a normal distribution unless your sample size is large, in which case you get an approximate normal distribution anyway by the Central Limit Theorem (remember those three words from Stats I?). Here, you aren’t making an assumption, but examining a *condition* (something you check before proceeding). You plot the data, see if they meet the condition, and if they do, you proceed. If not, you can use nonparametric methods instead (discussed in Chapter 16).



Nearly every statistical technique for analyzing data has at least some conditions on the data in order for you to use it. Always find out what those conditions are, and check to see whether your data meet them (and if not, consider using nonparametric statistics; see Chapter 16). Be aware that many statistics textbooks wrongly use the word *assumption* when they actually mean *condition*. It’s a subtle but very important

difference.

Only Reporting “Important” Results



As a data analyst, you must not only avoid the pitfall of reporting only the significant, exciting, and meaningful results but also be able to detect when someone else is doing so. Some number crunchers examine every possible option and look at their data in every possible way before settling on the analysis that gets them the desired result.

You can probably see the problem with that approach. Every technique carries the chance for error. If you’re doing a *t*-test, for example, and the α level is 0.05, over the long term 5 out every 100 *t*-tests you conduct will result in a false alarm (meaning you declare a statistically significant result when it wasn’t really there) just by chance. So, if an eager researcher conducts 20 hypothesis tests on the same data set, odds are that at least one of those tests could result in a false alarm just by chance, on average. As this researcher conducts more and more tests, he’s unfairly increasing his odds of finding something that occurred just by chance and running the risk of a wrong conclusion in the process.

It’s not all the eager researcher’s fault, though. He’s pressured by a results-driven system. It’s a sad state of affairs when the only results that get broadcast on the news and appear in journal articles are the ones that show a statistically significant result (when H_0 is rejected). Perhaps it was a bad move when statisticians came up with the term *significance* to denote rejecting H_0 — as if to say that rejecting H_0 is the only important conclusion you can come to. What about all the times when H_0 couldn’t be rejected? For example, when doctors failed to conclude that drinking diet cola causes weight gain, or when pollsters didn’t find that people were unhappy with the president? The public would be better served if researchers and the media were encouraged to report the statistically insignificant but still important results along with the statistically significant ones.



The bottom line is this: In order to find out whether a statistical conclusion is correct, you can’t just look at the analysis the researcher is showing you. You also have to find out about the analyses and results they’re *not* showing you and ask questions. Avoid the urge to rush to reject H_0 .

Assuming a Bigger Sample Is Always Better

Bigger is better in some things, but not always with sample sizes. On one hand, the bigger your sample is, the more precise the results are (if no bias is present). A bigger

sample also increases the ability of your data analysis to detect differences from a model or to deny some claim about a population (in other words, to reject H_0 when you're supposed to). This ability to detect true differences from H_0 is called the *power* of a test (see Chapter 3). However, some researchers can (and often do) take the idea of power too far. They increase the sample size to the point that even the tiniest difference from H_0 sends them screaming to press that all-important reject H_0 button.



Sample sizes should be large enough to provide precision and repeatability of your results, but there's such a thing as being too large, believe it or not. You can always take sample sizes big enough to reject any null hypothesis, even when the actual deviation from it is embarrassingly small. What can you do about this? When you read or hear that a result was deemed statistically significant, ask what the sample mean actually was (before it was put into the t -formula) and judge its significance to you from a practical standpoint. Beware of someone who says, "These results are statistically significant, and the large sample size of 100,000 gives even stronger evidence for that."

Suppose research claims that the typical in-house dog watches an average of ten hours of TV per week. Bob thinks the true average is more, based on the fact that his dog Fido watches at least ten hours of cooking shows alone each week. Bob sets up the following hypothesis test: $H_0: \mu = 10$ versus $H_a: \mu > 10$. He takes a random sample of 100 dogs and has their owners record how much TV their dogs watch per week. The result turns out a sample mean of 10.1 hours, and the sample standard deviation is 0.8 hours. This result isn't what Bob hoped for because 10.1 is so close to 10. He calculates the test statistic for

$$t = \frac{(\bar{x} - \mu)}{\frac{s}{\sqrt{n}}} = \frac{(10.1 - 10.0)}{\frac{0.8}{\sqrt{100}}} = \frac{0.1}{0.08}$$

this test using the formula and comes up with a value of , which equals 1.25 for t . Because the test is a right-tailed test ($>$ in H_a), Bob can reject H_0 at α if t is beyond 1.645, and his t -value of 1.25 is far short of that value. Note that because $n = 100$ here, you find the value of 1.645 by looking at the very last row of the t -distribution table (visit Table A-1 in the appendix). The row is marked with the infinity sign to indicate a large sample. So Bob can't reject H_0 . To add insult to injury, Bob's friend Joe conducts the same study and gets the same sample mean and standard deviation as Bob did, but Joe uses a random sample of 500 dogs rather than 100. Consequently, Joe's t -value is

$$t = \frac{(10.1 - 10.0)}{\frac{0.8}{\sqrt{500}}} = \frac{0.1}{0.036}$$

, which equals 2.78. Because 2.78 is greater than 1.645, Joe gets to reject H_0 (to Bob's dismay).



Why did Joe's test find a result that Bob's didn't? The only difference was the sample size. Joe's sample was bigger, and a bigger sample size always makes the standard error smaller (see Chapter 3). The standard error sits in the denominator of

the t-formula, so as it gets smaller, the t-value gets larger. A larger t-value makes it easier to reject H_0 . (See Chapter 3 for more on precisions and margin of error.)

Now, Joe could technically give a big press conference or write an article on his results (his mom would be so proud), but you know better. You know that Joe's results are technically statistically significant, but not practically significant — they don't mean squat to any person or dog. After all, who cares that he was able to show evidence that dogs watch just a tiny bit more than ten hours of TV per week versus exactly ten hours per week? This news isn't exactly earth-shattering.

It's Not Technically Random, But . . .

When you take a sample on which to build statistical results, the operative word is *random*. You want the sample to be randomly selected from the population. A problem is that people often collect a sample that they think is mostly random, sort of random, or random enough — and that doesn't cut it. A plan for taking a sample is either random or it isn't.

One day I gave each of the 50 students in my class a number from 1 to 50, and I drew two numbers randomly from a hat. The two students I picked were sitting in the first row, and not only that, they were right next to each other. My students immediately cried foul!

After this seemingly odd result, I took the opportunity to talk to my class about truly random samples. A *random sample* is chosen in such a way that every member of the original population has an equal chance of being selected. Sometimes people who sit next to each other are chosen. In fact, if these seemingly strange results never happen, you may worry about the process; in a truly random process, you're going to get results that may seem odd, weird, or even fixed. That's part of the game.

In my consulting experiences, I always ask how my clients chose or plan to choose their samples. They always say they'll make sure it's random. But when I ask them how they'll do this, I sometimes get less-than-stellar answers. For example, someone needed to get a random sample from a population of 500 free-range chickens in a farmyard. He needed five chickens and said that he'd select them randomly by choosing the five that came up to him first. The problem is, animals that come up to you may be friendlier, more docile, older, or perhaps more tame. These characteristics aren't present in every chicken in the yard, so choosing a sample this way isn't random. The results are likely biased in this case.



Always ask the researcher how he or she selected a sample, and when you select your own samples, stay true to the definition of random. Don't use your own judgment to choose a random sample; use a computer to do it for you!

Assuming That 1,000 Responses Is 1,000 Responses

A newspaper article on the latest survey says that 50 percent of the respondents said blah blah blah. The fine print says the results are based on a survey of 1,000 adults in the United States. But wait — is 1,000 the actual number of people selected for the sample, or is it the final number of respondents? You may need to take a second look; those two numbers hardly ever match.

For example, Jenny wants to know the percentage of Americans who have ever knowingly cheated on their taxes. In her statistics class, she found out that if she gets a sample of 1,000 people, the margin of error for her survey is only ± 3 percent, which she thinks is groovy. So she sets out to achieve the goal of 1,000 responses to her survey. She knows that in these days it's hard to get people to respond to a survey, and she's worried that she may lose a great deal of her sample that way, so she has an idea. Why not send out more surveys than she needs, so that she gets 1,000 surveys back?

Jenny looks at several survey results in the newspapers, magazines, and on the Internet, and she finds that the response rate (the percentage of people who actually responded to the surveys) is typically around 25 percent. (In terms of the real world, I'm being generous with this number, believe it or not. But think about it: How many surveys have you thrown away lately? Don't worry, I'm guilty of it too.) So, Jenny does the math and figures that if she sends out 4,000 surveys and gets 25 percent of them back, she has the 1,000 surveys she needs to do her analysis, answer her question, and have that small margin of error of ± 3 percent.

Jenny conducts her survey, and just like clockwork, out of the 4,000 surveys she sends out, 1,000 come back. She goes ahead with her analysis and finds that 400 of those people reported cheating on their taxes (40 percent). She adds her margin of error, and reports, "Based on my survey data, 40 percent of Americans cheat on their taxes, ± 3 percentage points."

Now hold the phone, Jenny. She only knows what those 1,000 people who returned the survey said. She has no idea what the other 3,000 people said. And here's the kicker: Whether or not someone responds to a survey is often related to the reason the survey is being done. It's not a random thing. Those nonrespondents (people who don't respond to a survey) carry a lot of weight in terms of what they're not taking time to tell you.

For the sake of argument, suppose that 2,000 of the people who originally received the survey were uncomfortable with the question because they *do* cheat on their taxes; they just didn't want anyone to know about it, so they threw the survey in the trash. Suppose that the other 1,000 people don't cheat on their taxes, so they didn't think it was an issue and didn't return the survey. If these two scenarios were true, the results would look like this:

$$\text{Cheaters} = 400 \text{ (respondents)} + 2,000 \text{ (nonrespondents)} = 2,400$$

These results raise the total percentage of cheaters to 2,400 divided by 4,000 = 60 percent. That's a huge difference!

You could go completely the other way with the 3,000 nonrespondents. You can suppose that none of them cheat, but they just didn't take time to say so. If you knew this info, you would get 600 (respondents) + 3,000 (nonrespondents) = 3,600 noncheaters. Out of 4,000 surveyed, this would mean 90 percent didn't cheat, and only 10 percent did. The truth is likely to be somewhere between the two examples I just gave you, but nonrespondents make it too hard to tell.

And the worst part is that the formulas Jenny uses for margin of error don't know that the information she put into them is based on biased data, so her reported 3 percent margin of error is wrong. The formulas happily crank out results no matter what. It's up to you to make sure that you put good, clean info into the formulas.



Getting 1,000 results when you send out 4,000 surveys is nowhere near as good as getting 1,000 results when sending out 1,000 surveys (or even 100 results from 100 surveys). Plan your survey based on how much follow-up you can do with people to get the job done, and if it takes a smaller sample size, so be it. At least the results have a better chance of being statistically on target.

Of Course the Results Apply to the General Population

Making conclusions about a much broader population than your sample actually represents is one of the biggest no-no's in statistics. This kind of problem is called *generalization*, and it occurs more often than you may think. People want their results instantly; they don't want to wait for them, so well-planned surveys and experiments take a back seat to instant Web surveys and convenience samples.

For example, a researcher wants to know how cable news channels have influenced the way Americans get their news. He also happens to be a statistics professor at a large research institution and has 1,000 students in his class. He decides that instead of taking a random sample of Americans, which would be difficult, time-consuming, and expensive, he'll just put a question on his final exam to get his students' answers. His data analysis shows him that only 5 percent of his students read the newspaper and/or watch network news programs; the rest watch cable news. For his class, the ratio of students who exclusively watch cable news compared to those students who don't is 20 to 1. The professor reports this and sends out a press release about it. The cable news channels pick up on it and the next day are reporting, "Americans choose cable news channels

over newspapers and network news by a 20 to 1 margin!"

Do you see what's wrong with this picture? The professor's conclusions go way beyond his study, which is wrong. He used the students in his statistics class to obtain the data that serves as the basis for his entire report and the resulting headline. Yet the professor reports that these results are true for all Americans. I think it's safe to say that a sample of 1,000 college students taking a statistics class at the same time at the same college doesn't represent a cross section of America.

If the professor wants to make conclusions in the end about America, he has to select a random sample of Americans to take his survey. If he uses 1,000 students from his class, his conclusions can only be made about that class and no one else.



To avoid or detect generalization, identify the population that you're intending to make conclusions about and make sure the sample you selected represents that population. If the sample represents a smaller group within that population, you have to downsize the scope of your conclusions also.

Deciding Just to Leave It Out

It seems easier sometimes to just leave out information. I see this all too often when I read articles and reports based on statistics. But, this error isn't the fault of only one person or group. The guilty parties can include

✓ **The producers:** Some researchers may leave statistical details out of their reports for a variety of reasons, including time and space constraints. After all, you can't write about every element of the experiment from beginning to end. However, other items they leave out may be indicative of a bigger problem. For example, reports often say very little about how they collected the data or chose the sample. Or they may discuss the results of a survey but not show the actual questions they asked. Ten out of 100 people may have dropped out of an experiment, and the researchers don't tell you why. All these items are important to know before making a decision about the credibility of someone's results.

Another way in which some data analysts leave information out is by removing data that doesn't fit the intended model (in other words, "fudging" the data). Suppose a researcher records the amount of time spent surfing the Internet and relates it to age. He fits a nice line to his data indicating that younger people surf the Internet much more than older people and that surf time decreases as age increases. All is good except for Claude the outlier who's 80 years old and surfs the Internet day and night, leading his own bingo chat rooms and everything. What to do with Claude? If not for him, the relationship looks beautiful on the graph; what harm would it do to remove him? After all, he's only one person, right?



No way. Everything is wrong with this idea. Removing undesired data points from a data set is not only very wrong but also very risky. The only time it's okay to remove an observation from a data set is if you're certain beyond doubt that the observation is just plain wrong. For example, someone writes on a survey that she spends 30 hours a day surfing the Internet or that her IQ is 2,200.

- ✓ **The communicators:** When reporting statistical results, the media leave out important information all the time, which is often due to space limitations and tight deadlines. However, part of it is a result of the current, fast-paced society that feeds itself on sound bites. The best example is survey results in which the margin of error isn't communicated. You can't judge the precision of the results without it.
- ✓ **The consumers:** The general public also plays a role in the leave-things-out mindset. People hear a news story and instantly believe it to be true, ignoring any chance for error or bias in the results. For example, you need to make a decision about what car to buy, and you ask your neighbors and friends rather than examine the research and the resulting meticulous, comprehensive ratings. At one time or another everyone neglects to ask questions as much as they should, which indirectly feeds the entire problem.

In the chain of statistical information, the producers (researchers) need to be comprehensive and forthcoming about the process they conducted and the results they got. The communicators of that information (the media) need to critically evaluate the accuracy of the information they're getting and report it fairly. The consumers of statistical information (the rest of us) need to stop taking results for granted and to rely on credible sources of statistical studies and analyses to help make important life decisions.



In the end, if a data set looks too good, it probably is. If the model fits too perfectly, be suspicious. If it fits exactly right, run and don't look back! Sometimes what's left out speaks much louder than what's put in.

Chapter 22

Ten Ways to Get Ahead by Knowing Statistics

In This Chapter

- ▶ Knowing what information to look for
 - ▶ Being skeptical and confident
 - ▶ Piecing together the statistics puzzle and checking your answers
 - ▶ Knowing how best to present your findings
-

One of my personal goals of teaching statistics is to help people get very good at being able to say “Wait a minute!” and stop a wrong analysis or a misleading graph in its tracks. I also want to help them become the stats gurus in their workplaces — those people who aren’t afraid to work with statistics and do so correctly and confidently (and to also know when to consult a professional statistician). This chapter arms you with ten ways of trusting your statistics instincts and increasing your professional value through your understanding of the critical world of stats.

Asking the Right Questions

Every study, every experiment, and every survey is done because someone had a question they wanted answered. For example “How long should this warranty last?”; “What’s the chance of me developing complications during surgery?”; “What does the American public think about banning public smoking?” Only after a clear question has been defined can proper data collection begin.

Suppose a restaurant owner tells me that he wants to conduct a survey to learn more about the clientele at his restaurant. We talk about various variables to look at, including the number of people in the party, how often they’ve been there before, the type of food ordered, the amount they pay, how long they stay, and so on. After we collect some data and go over the results, he suddenly has a major realization: What he really wants to do is compare the clientele of his lunch crowd to his dinner crowd. Does the dinner crowd spend more money? Are they older? Do they stay longer? But sadly, he can’t answer any of those questions because he didn’t mention collecting data on whether the customers were there for lunch or dinner.

What happened here is a common mistake. The restaurant owner said he “just wanted to study” his clientele; he never mentioned comparisons because he hadn’t thought that far ahead. If he had thought about it, he would have realized the real question was, “How does my lunch clientele compare with my dinner clientele?” Then including a question on whether diners were there for lunch or dinner would have been a no-brainer. Always ask the right questions to get the answers you need.



Testing the waters a bit before plunging into a full-blown study can be very helpful. One way to do this is to conduct what researchers call a pilot study. A pilot study is a small exploratory study that you use as a testing ground for the real thing. For example, you design a survey and try it out on a small group to see if they find any confusing questions, redundancies, spelling errors, and so on. Pilot studies are a quick and inexpensive way to help ensure that all goes well when the actual study takes place.

Being Skeptical

Being statistically skeptical is a good thing (within reason). Some folks have given up on statistics, thinking that people can say anything they want if they manipulate the data enough. So those who have a healthy degree of skepticism can get ahead of the game.

Colorful charts and graphs can catch your eye, especially if they have neat little captions, and long and detailed professional reports may show you more information than you want to know, all laid out in neat tables, page after page. What's most important, however, is not how nice-looking the information is, or how professionally sound or scientific it looks. What's most important is what's happened behind the scenes, statistically speaking, in order to produce results that are correct, fair, and clear.



Many folks know only enough statistics to be dangerous. And many reported results are incorrect, either by mistake or by design (unfortunately). It's better to be skeptical than sorry!

Here's how to put your skepticism to good use:

- ✓ Get a copy of survey questions asked. If the questions are misleading, the survey results aren't credible.
- ✓ Find out about the data-collection process. When was the survey conducted? Who was selected to participate? How was the information collected? Surveys conducted on the Internet and those based on call-in polls are almost always biased, and their results should be thrown out the window.
- ✓ Find out about the response rate of the survey. How many people were initially contacted? How many responded? If many were contacted and few responded, the results are almost certainly biased because survey respondents typically have stronger feelings than those who don't respond.

Collecting and Analyzing Data Correctly

On one hand it's very important to think very critically and even be skeptical at times about statistical results that you come across in everyday life and in the workplace. You should always ask questions before you deem the results to be credible.

On the other hand, it's also very important to remember that others are thinking critically about your results also, and you need to avoid the skepticism that you see others receiving. To avoid potential potshots that may be taken at your results, you need to make sure you've done everything right.

Because you're reading this book, by now you should have many tools to help you do data collection and analysis correctly. In each chapter you hear the same theme song: Using the wrong analysis or too many analyses isn't good. For each type of analysis I present, you see how to check to make sure that particular analysis is okay to use with the data you have. Chapters 1 and 2 serve as a reference to which techniques are needed, and where to find them in the book.



Ninety percent of the work involved in a statistical analysis happens before the data even goes into the computer. Here's a basic to-do list of what to check for:

- ✓ Design your survey, your experiment, or your study to avoid bias and ensure precision.
- ✓ Make sure you conduct the study at the right time and select a truly random sample of individuals to participate.
- ✓ Follow through with those participants to make sure your final results have a high response rate.

This to-do list can be challenging, but in the end, you'll be safe in knowing that your results will stand up to criticism because you did everything right.

Calling for Help

One of the toughest things for nonstatisticians to get is that they don't have to do all the statistics themselves. In fact, it's not a good way to go in many instances. The six most important words for any nonstatistician are "Know when to consult a statistician." Know when to ask for help. And the best time to ask for help is *before* you collect any data.

So how can you tell when you're in a bit over your head and you need someone to throw you a statistical lifeline? Here are some examples to help give you an idea of when to call:

- ✓ If your boss wants no less than a 100-page marketing results report on her desk by Monday and you haven't collected data point #1, CALL.
- ✓ If you're reading *Cosmopolitan* on your lunch break and you want to analyze how

you and your friends came out on the “Who’s the Gossip Queen in Your Workplace?” quiz, DON’T CALL.

- ✓ If the list of questions on your survey becomes longer than you are tall, CALL.
- ✓ If you want to make a bar graph of how many of your Facebook friends are fans of the 1970s, 80s or 90s, DON’T CALL.
- ✓ If a scatterplot of your data looks like it should be in a Rorschach inkblot test, CALL (and fast!).
- ✓ If you want to know the odds that someone you haven’t seen since high school is on the same plane to Africa as you are, DON’T CALL.
- ✓ If you have an important job to do that involves statistics and you are unsure of how to begin or how you’ll analyze your data once you get it, CALL. The sooner you call, the more the professionals can do to help you look good!

Retracing Someone Else’s Steps

At some point in your work life, you’ll take out a report, read it, and you’ll have a question about it. You’ll go to find the data, and after much searching, you’ll bring up a spreadsheet with rows upon rows and columns upon columns of numbers and characters. Your eyes will glaze over; you’ll have no idea what you’re looking at. You’ll tell yourself not to panic and just to find the person who entered all the data and find out what’s going on.

But then comes the bad news. Someone named Bob collected the data and entered it a couple of years ago, and Bob doesn’t work for the company anymore. Now what do you do? More than likely, you’ll have to ditch the data and the report, start all over again from scratch, and lose valuable time and money in the process.

How could this disaster have been prevented? All the following issues should have been addressed before Bob passed on his report:

- ✓ The report should include a couple of paragraphs telling how and when the data were collected, the names of the variables in the data set, where they’re located in the spreadsheet, and what their labels are.
- ✓ The report should include a note about missing data. Missing data are sometimes left blank, but they also can be written as a negative sign (-) or a decimal point. (Using zeroes for missing data is a special no-no because they will be confused with actual data values that equal zero.)
- ✓ The rows of the data set should be defined. For example, does each row represent one person? Do they have ID numbers?



Unfortunately, many people create statistical reports and then disappear without a trace, leaving behind a data mess that often can't be fixed. It's common courtesy to take steps to avoid leaving other people in the lurch, the way Bob did. Always leave a trail for the next person to pick up where you left off. And on the flipside, always ask for the explanation and background of a data set before using it.

Putting the Pieces Together

You should never jump right into an analysis expecting to get a one-number answer and then walk away. Statistics requires much more work! You should view every statistical problem as a puzzle whose pieces need to be put together before you can see the big picture of what's really going on.

For example, suppose a coffee vendor wants to predict how much coffee she should have ready for an upcoming football game in Buffalo, New York. Her first step is to think about what variables may be related to coffee sales. Variables may be cost of the coffee, ease of carrying it, seat location (who wants to walk a mile for a cup of coffee?), and age of fans. The vendor also suspects that temperature at the game may affect coffee sales, with low temps translating into higher sales.

The vendor collects data on all these variables and explores the relationships. She finds that coffee sales and temperature are somewhat related. But is there more to this story than temperature?

To find out, the vendor compares coffee sales for two games with the same temperature and notices a big difference. Looking deeper, she notices one game was on a Sunday and one was on a Monday. Attendance was higher on Monday, and that game had more adults in attendance. By analyzing the data, the vendor found that temperature is related to coffee sales, but so is attendance, day of the week that the game was played, and age of the fans. Knowing this information, the vendor was able to predict coffee sales more accurately with a lower chance of running out of coffee or wasting it. This example illustrates that putting the pieces together to keep an eye on the big picture can really pay off.

Checking Your Answers

After your data have been analyzed and you get your results, you need to take one more step before running giddily to your boss saying, "Look at this!" You have to be sure that you have the right answers.



By right answers, I don't mean that you need to have the results that your boss wants to hear (although that would be great, of course). Rather, you need to make sure your data analysis and calculations are correct and don't leave you high and dry when the questions start to come. Follow these basic steps:

1. Double-check that you entered the data correctly, and weed out numbers that obviously make no sense (such as someone saying that he's 200 years old, or that he sold 500 billion light bulbs at his store last year).

Mistakes influence the data and the results, so catch them before it's too late.

2. Make sure that your numbers add up when they're supposed to.

For example, if you collected data on number of employees for 100 companies and you don't list enough number groups to cover them all, you're in trouble! Also be on the lookout for data on an individual that have been entered twice. This error shows up if you sort the data by rows.

3. If you intend to make conclusions, make sure you're using the right numbers to do so.

If you want to talk about how crime has increased in your area over the last five years, showing the number of crimes on a graph is incorrect. The number of crimes can increase simply because the population size increases. For correct statistical conclusions about crime, you need to report the crime rate, which is the number of crimes per person (per capita), or the number of crimes per 100,000 people. Just take the number of crimes divided by the population size, or divided by 100,000, respectively. This approach takes population size out of it.

Explaining the Output

Computers certainly play a major role in the process of collecting and analyzing data. Many different statistical software packages exist, including MS Excel, Minitab, SAS, SPSS, and a host of others. Each type has its own style of printing out results. Understanding how to read, interpret, and explain computer output is an art and a science that not everyone possesses. With your statistical knowledge, though, you can be that person!

Computer output is the raw form of the results of doing any statistical summary or analysis. It can be graphs, charts, scatterplots, tables, regression analysis results, an analysis of variance table, or a set of descriptive statistics. Often the analysis is labeled by the computer; for example, ANOVA indicates an Analysis of Variance has been conducted (see Chapter 9). Graphs, charts, and tables, however, require the user to tell the computer what labels, titles, or legends (if any) to include so that the audience can quickly understand what's what.

Interpreting computer output involves sifting through what can seem like an intimidating amount of information. The trick is knowing exactly what results you want and where the

computer places them on the output. For example, in the output from a regression analysis, you find the equation of the regression line by looking in the COEF column of the output (see Chapter 4).

Most of the time there's information on a computer output that you don't need; sometimes there's also information that you don't understand. Before skipping everything, you may want to consult a statistician to make sure you aren't missing an important step, such as examining the correlation coefficient before doing a regression analysis (see Chapter 4).



Explaining what you found on computer output involves sizing up the knowledge of the person you're talking to, too. If you're writing an executive report, you don't need to explain every little thing you did and why; just use the parts of the output that tell the bottom line, and explain how it affects the company. If you're helping a colleague understand results, give him some reference information about the analyses (you can use this book). For example, you may discuss what a histogram does in general before you talk about the results of your histogram.

Most importantly, make sure the analysis is correct before explaining it to anyone. Sometimes it's easy when analyzing data to click on the wrong variable or to highlight the wrong column of data, which makes the analysis totally wrong.

Making Convincing Recommendations

As one moves up the corporate ladder, the less time she has to read reports and carefully examine statistics. The best data analysis in the world won't mean squat if you can't communicate your results to someone who doesn't have the time or interest to get into the nitty-gritty. In this data-driven world, statistics can play a major role in good decision-making. The ability to use statistics to make an effective argument, make a strong case, or give solid recommendations is critical.

Put yourself in the following situation. You've done the work, you've collected marketing and sales data, and you've done the analyses and processed the results. Based on your study of product placement for your Sugar Surge Pop, you determined that the best strategy for placing this product on grocery store shelves is to put it in the checkout aisle at eye level so children can see it. (You never see nail clippers or hand sanitizers on the kid's eye-level shelves in the impulse aisle, do you?) Word is that your boss favors putting this product in the candy aisle of the store. (Of course she has no data to support this, just her own experience people-watching in the candy aisle.) How do you convince her to follow your recommendation?

Probably the worst thing you can do is go into her office with a 100-page report loaded with everything from soup to nuts. Loads of complex information may impress your

mom, but it won't impress your boss. Save that report in case she asks for it (or in case you need a doorstop). What you need is a short, succinct, and straightforward report that makes the point. Here's how to craft it:

1. Start out with a statement of the problem.

"We want to determine which location has the most sales of the Sugar Surge Pop."

2. Briefly outline your data-collection process.

"We chose 50 stores at random and placed the product in the checkout aisle at 25 stores and in the candy aisle in the other 25. We controlled for other factors such as number of products placed."

3. Describe what data you collected.

"We tracked the sales of the product over a six-month period, calculating weekly sales totals for each store." At this point, show your charts and graphs of the sales over time for the two groups.

4. Tell briefly how you analyzed the data, but spend the most time on your findings.

Don't show the output — your boss doesn't need to see that. You know the expression, "Never let them see you sweat"? That's important here. What you do want to say is, "We did a statistical analysis comparing average sales at these locations, and we found sales in the checkout aisle to be significantly higher than sales in the candy aisle." You can quantify the difference with percentages.

Follow up with your recommendation for product placement at kids' eye-level in the checkout aisle, being sure to answer the original question you started with in Step 1. Then the most important tip is to let your boss think the optimal placement was all her idea!

Establishing Yourself as the Statistics Go-To Guy or Gal

Nothing is more valuable than someone in the workplace who isn't afraid to do statistics. Every office has one person with the courage to calculate, the confidence to make confidence intervals, the willingness to wrestle with the output, and the gumption to graph. This person is eventually everyone's friend and the first person to get to know when starting a new job.

What are the perks of being the statistics go-to guy or gal? It's the glory of knowing that you're saving the day, taking one for the team, and standing tall in the face of disaster. Your colleagues will say, "I owe you one," and you can take them up on it.

But seriously, the statistics go-to guy has a more secure job because his boss knows that statistics is a staple of the workplace, and having someone to jump in when needed is invaluable.



Statistics and statistical analyses can be intimidating, yet they're critical for the workplace. In most any career these days, you need to know how to select samples, write surveys, set up a process for collecting the data, and analyze it.

Chapter 23

Ten Cool Jobs That Use Statistics

In This Chapter

- ▶ Using statistics in a wide range of jobs (yours might be next!)
- ▶ Tracking data from birds to sports to crime
- ▶ Taking stats into the professional world of medicine, law, and the stock market

This book is meant to be a guide for folks who need to know statistics for their everyday life (which is all of us) as well as in their workplaces (which is most of us). If I think about it long enough, I can come up with some use of statistics in almost every job out there (except maybe a psychic advisor).

This chapter features a cross section of ten careers that all involve statistics in some way, shape, or form. You may be surprised at how often statistics turn up in the workplace! So don't burn this book when your stats course is over; you may find it to be useful in your job hunt or your job. (My accountant has a copy of this book on his shelf — what does that say? As long as he doesn't have a copy of *Accounting For Dummies* alongside it, I guess we'll be okay.)

One of my personal goals as a teacher of statistics is to help my students become the go-to folks in the workplace. You know, that person with a background in statistics who knows what she's doing, and when it's crunch time, she can do the statistics needed correctly and confidently. With experience and help from this book, you too can become that person. You'll become a hero, and your job will be all the more secure for it.

Pollster

Pollsters collect information on people from populations they're interested in. Some of the big names in professional polling include the Gallup Organization, the Associated Press (AP), Zogby International, Harris Interactive, and the Pew Research Center. Major news organizations such as NBC, CBS, and CNN also conduct polls, as do many other agencies and organizations.

The purposes of polls vary from the medical field trying to determine what's causing obesity, to political pollsters who want to keep up with the daily pulse of American opinion, to surveys that provide feedback and ideas to corporations.



Knowledge of statistics is considered golden in the polling industry, because jobs can include designing surveys; selecting a proper sample of participants; carrying

out a survey to collect data; and then recording, analyzing, and presenting the results.

All these tasks are part of statistics — the art and science of collecting and making sense of data. But don't just take my word for it; here's a quote from a job posting for the Gallup Organization for a Research Analyst. I have to say it totally screams STATISTICS!

If you have a strong academic record in the social sciences or economics, a familiarity with quantitative and categorical research and statistical tools in Market Research/Survey Research or Consulting, enjoy pulling together research data and abstract concepts to tell a meaningful story, while continually learning — this is the place to manage processes and projects that deliver perfect completion of client engagements.

And here's something you don't see every day. I found a job posting for a polling analyst with roughly the same requirements but a very different work setting. The job was for a company that provides security and intelligence for the United States government. For this job you need federal security clearance. You never know where your statistical background is going to take you!

Other positions related to polling that I've seen listed are quantitative research specialist and public polling research analyst.



A great Web site for finding out more about what pollsters do and what their work looks like is, appropriately, www.pollster.com.

Ornithologist (Bird Watcher)

Everyone watches birds on occasion. I gladly admit to being a semi-serious birdwatcher, always trekking to Magee Marsh on Lake Erie in May for International Migratory Bird Day. But have you ever thought about getting paid to watch birds and other wildlife? Today's ever-increasing awareness of the environment includes a great deal of focus on identifying, studying, and protecting wildlife of all kinds.

Ornithology is the science of bird study. Ornithologists are always collecting data and finding and studying statistics on birds — often on a certain type of bird and its behavior. Some examples of common bird statistics include:

- ✓ Bird counts (number of birds per square unit of space on a particular day)
- ✓ Nest locations and territory maps
- ✓ Number of eggs laid and hatched
- ✓ Food preferences and foraging techniques

- ✓ Behaviors caught on tape and quantified



You can tap into a Web site totally dedicated to jobs that need birdwatchers and wildlife watchers and the use of their statistical skills. Here's one of the job postings supplied by the Ornithological Society:

FLAMMULATED OWL SURVEY TECHNICIANS (2) needed for Idaho Bird Observatory study of Flammulated Owls and other forest birds in Idaho (approx. 2.5 months). Duties will consist mainly of standardized surveys and data entry. Qualifications of applicants should include: 1) good eyesight and hearing, 2) proficiency with standardized survey procedures, 3) ability to identify Western birds by sight and sound, and 4) willingness to give your all. Candidates must be physically fit and undaunted by the prospects of heat, humidity, bugs, and mud. (Indeed!)

With more experience and knowledge, you can eventually become a research wildlife biologist for the U.S. Department of the Interior U.S. Geological Survey. A job description for this position just found today actually requires 15 credit hours of statistics, proving that the government is onto the whole statistics thing.

Sportscaster or Sportswriter

Every good sportscaster or sportswriter knows that you're nothing without good juicy statistics that no one else knows. You do your homework by studying training camps and pouring over printouts, spreadsheets, and historical data. You read newspapers, look at record books, and watch film. There's no shortage of data out there, and your audience can't get enough.

Sports fans are statistics addicts! (Being an Ohio State Buckeye, I'm as rabid as the rest of 'em.) Here's just a sampling of the statistics recorded and presented in my favorite sport of college football:

- ✓ Points scored
- ✓ Points against
- ✓ Rushing yards
- ✓ Receiving yards
- ✓ Passing yards
- ✓ Interceptions
- ✓ Fumbles
- ✓ Punt and kick returns

- ✓ Number and distance of field goals attempted and made
- ✓ Kicker's career longest
- ✓ Number of first downs
- ✓ Third down conversions
- ✓ Fourth down conversions
- ✓ Penalties
- ✓ All-purpose yards
- ✓ Total offense
- ✓ Total defense
- ✓ Number of sacks
- ✓ Rushing defense
- ✓ Passing defense
- ✓ Turnover margins
- ✓ Passing efficiency
- ✓ Scoring offense
- ✓ Scoring defense
- ✓ Scoring by special teams
- ✓ Coaches' Poll standings
- ✓ AP Poll standings
- ✓ BCS standings (that's another book for another day!)
- ✓ The most 12+ win seasons
- ✓ Coaching records
- ✓ Single game high scores
- ✓ Game attendance
- ✓ Toughness of schedule
- ✓ Winning and losing streaks
- ✓ Coin toss winners

It's obvious that we need a new saying: "Those who play sports play. Those who watch sports do statistics."

Journalist

Journalists of all types at some point or another have to work with data. They have a hard row to hoe, because the data comes to them from an infinite number of possible avenues and channels on an infinite number of topics, and they need to make sense of it, pick out what they feel is most important, boil it down, present the results, and write a story around it, all under a very strict deadline (sometimes only a few hours). That's a big job!

As a consumer of the media, I see many good uses of statistics that are clear, correct, and make interesting and important points. However, I also see many incorrect and misleading statistics in the media, and I cringe every time.

Some of the most common problems include making simple math errors, reporting percentages above 100 percent, assuming cause and effect relationships that aren't proven, using misleading graphs, and leaving out information (such as the number of people surveyed, the rate of nonresponse, and the margin of error). But for me the biggest problem is reading a headline that sounds catchy, eye-opening, perhaps even shocking, only to find that it's not corroborated by the statistics in the article.

Having a couple of solid statistics courses under your belt puts you way ahead in that job interview for a journalist position. The statisticians around the world are counting on you to get out there on your white horse and do things right! (Don't forget to take this book with you in your saddle bag!)

To recognize the importance and appreciation of the difficult task that journalists have in using and reporting with statistics, the Royal Statistical Society has established an Award for Statistical Excellence in Journalism. Following is the description for the award, and I couldn't agree with it more!

The Royal Statistical Society wishes to encourage excellence in journalists' use of statistics to question, analyze, and investigate the issues that affect society at large. Journalistic excellence in statistics helps to hold decision makers in all sectors to account — through accessible communication of complex information, highlighting of success, and exposure of important missing information.

Crime Fighter

Crime statistics help the nation's crime fighters, such as police officers, determine which kinds of crimes occur where, how often, to whom, and by whom. Crime statistics for the entire nation are compiled and analyzed by the U.S. Department of Justice. National Crime Victimization Surveys are also conducted to help understand trends in crimes of various types.

Police officers record every incident they're involved in, forming large databases that

city, county, and state officials can use to determine the number of police officers needed and which areas to focus most heavily on, and also to make changes in the policies and procedures of their police departments. The FBI can also use these huge databases to track criminals, look for patterns in crime types and occurrences, and keep track of overall trends in the number of crimes as well as the type of crimes that occur over time.

People looking for a new home or a new school can consult freely available information on crime statistics, and politicians use it to show that crime is going up or down, that money should or should not be spent on more police officers, and how safe their city or state has become with them in office.

Here's an overview of how the U.S. Department of Justice Web site uses data and statistics to help fight crime:

All states have established a criminal record repository which maintains criminal records and identification data and responds to law enforcement inquiries and inquiries for other purposes such as background checks and national security. Criminal records include data provided by all components of the criminal justice system: law enforcement, prosecution, courts, and corrections. . . . Records developed for statistical purposes describe and classify each criminal incident and include data on offender characteristics, relationships between the offender and the victim, and offense impact. Statistical data are extracted from operational records using uniform criteria for classification and collection. Detailed statistical data permit localities to identify problem areas and to allocate manpower and limited financial resources in an efficient and effective manner.

Medical Professional

People who work in the medical field depend on statistics to do research and find new cures, therapies, medicines, and procedures to increase the health and well-being of all people. Medical researchers conduct clinical trials to measure every conceivable side effect of every drug that goes through the approval process. Comparative studies are done all the time to determine what factors influence weight, height, intelligence level, and ability to survive a certain disease. Statistics are a lifeline to being more confident that what works for a sample of individuals will also work for the population for which it was meant.

In the medical profession, the use of statistics starts as soon as a patient's name is called in the waiting room. Suppose you're a nurse. The first thing you ask the patient to do is to "go ahead and step on the scale please" (eight dreaded words heard in doctors' offices round the world). From there, you check his vital signs (also known as vital statistics): temperature, blood pressure, pulse rate, and sometimes, respiration (breathing) rate. You record his numbers in your computer and compare them to what has been determined to be the normal range. Setting the normal ranges involves statistics as well, through analyzing historical data and medical research.

Like other kinds of statistics, the way in which a person's vital statistics are collected can greatly affect the results. For example, different types of blood pressure instruments exist, and some are better than others. And we all know the scale in the doctor's office is always at least 10 pounds over!

Marketing Executive

Marketing is critical to any product's success. That's why companies spend millions of dollars for 30-second commercials during the Super Bowl. Researching who will buy your product, where, when, and for how much is a job that includes lots of statistics.



Some data is what statisticians call *quantitative*, such as surveys of existing, past, and potential customers, sales information and trends, economic and demographic information, and data regarding competitors. Other data is *categorical*, including in-depth one-on-one interviews and focus groups to get a general picture of what consumers think, how they feel, what ideas they have, and what additional information they need about your product. (See Chapter 2 or your Stats I textbook for a review of quantitative and categorical data.)

Consider the example of Mars, Incorporated, which makes M&M'S candy. How has the company's product become a national icon for kids of all ages? The secrets have to be Mars' innate ability to change with the times and knowledge of what its customers want.

Statistics plays a huge role in this success through collecting data on sales, but most importantly, through getting direct feedback from customers using interviews, focus groups, and surveys. (I can't imagine the strain of trying out M&M'S and talking about what I think. "Oh wait, I need another sample before I can give you a good answer.") By analyzing this data, Mars is able to determine some of the most important and intricate details that spell success and longevity of any company.

For example, in 1995 Mars conducted a nationwide survey asking customers to choose the newest M&M'S color. That's when blue came on the scene. In 2002, the survey went global and purple became the new addition to the M&M'S palette. The Mars company uses statistics to find ways to be innovative in making new colors, flavors, styles, and even allowing for customized M&M'S, yet it still retains the classic essence of the M&M'S that started it all.

Lawyer

You've no doubt heard the phrase "beyond a reasonable doubt." It's the code that jurors use to make a decision of guilty or not guilty. The field of statistics plays a major role in determining whether laws are being followed or broken, whether a defendant is guilty or

innocent, and whether laws need to be created or changed. Statistical information is very powerful evidence.

Lawsuits are often settled on the basis of statistical evidence collected in multiple situations over the course of years. Statistics also allow lawmakers to break down information in order to propose new laws. For example, using statistics to show that the first two hours are the most critical in terms of finding a missing child led to the Amber Alerts broadcast on TV and radio and posted on highways when children go missing.

Prosecutors and defense attorneys often use probability and statistics to help make their cases, too. They also have to make these statistics understandable to a jury. (Maybe copies of this book should be a requirement for sequestered juries!) Statisticians are often brought onto the legal team to help attorneys gather information, decipher the results, and use the data in a jury trial situation.

Attorneys may use correlation to show that certain variables have a linear relationship, such as skid distance and amount of a certain type of concrete in pavement, or the strength of a bridge beam related to the weight placed upon it.

Statistics also can help test claims. For example, suppose Shipping Company A claims its packages are delivered on average two hours faster than Company B. If a random sample of packages takes longer than two days to arrive and the difference is large enough to have strong evidence against Company A's claim, it could get in trouble for false advertising. Of course, any decision based on statistics can be wrong, just by chance. In this case, if the random sample of packages just happened to take longer than usual and doesn't represent the typical average delivery time, Shipping Company A can fight back, saying they were unjustly accused of false advertising. It's a tight line to walk, and statisticians try their best to set up procedures to help the real truth come to light.

Stock Broker

Professional gamblers are folks who make a living by gambling. They know the ropes, they've been all over town, and they're good at what they do, which is basically play games that they feel they can win with skill as well as luck. There are different styles of professional gamblers. One type goes for casino games, such as poker and blackjack.

Another type of gambler is the stock broker who dresses up in a suit, has power lunches, and constantly changes his decisions throughout the day depending on the current state of the system. Stock brokers make predictions and decisions on buying and selling stocks on a minute-to-minute basis.

Statistics plays a critical role in a successful stock broker's decision making. To gain an edge, brokers use sophisticated data collection and analysis software, financial models, and numbers from all over the world, and they have to be able to analyze the information, interpret it, and act on it quickly.

However, one has to also keep an eye out for stock brokers who either ignore the real statistics or come up with their own statistics to give their clients an unrealistic view of what's going on with their money. There have even been instances where brokers have stolen their clients' money right out from under their noses to the tune of millions and even billions of dollars. These situations aren't common, but they make a huge impact on the confidence of investors and ultimately even on the state of the economy.

Collecting and using verifiable data, and analyzing that data with legitimate techniques, should be the goal of all good stock brokers (and for everyone who uses statistics in the workplace).

Appendix

Reference Tables

This appendix includes commonly used tables for five important distributions for Stats II: the *t*-distribution, the binomial distribution, the Chi-square distribution, the distribution for the rank sum test statistic, and the *F*-distribution.

t-Table

Table A-1 shows right-tail probabilities for the *t*-distribution (refer to Chapter 3). To use Table A-1, you need four pieces of information from the problem you're working on:

- ✓ The sample size, n
- ✓ The mean of x , denoted μ
- ✓ The standard deviation of your data, s
- ✓ The value of x for which you want the right-tail probability

After you have this information, transform your value of x to a *t*-statistic (or *t*-value) by taking your value of x , subtracting the mean, and dividing by the standard error (see

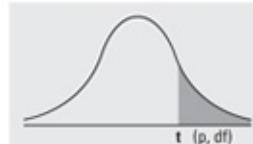
$$t_{n-1} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

Chapter 3) by using the formula .

Then look up this value of t on Table A-1 by finding the row corresponding to the degrees of freedom for the *t*-statistic ($n - 1$). Go across that row until you find two values between which your *t*-statistic falls. Then go to the top of those columns and find the probabilities there. The probability that t is beyond your value of x (the right-tail probability) is somewhere between these two probabilities. Note that the last row of the *t*-table shows $df = \infty$, which represents the values of the *Z*-distribution, because for large sample sizes t and *Z* are close.

Table A-1**The t-Table**

t-distribution showing area to the right



df/p	0.40	0.25	0.10	0.05	0.025	0.01	0.005	0.0005
1	0.324920	1.000000	3.077684	6.313752	12.70620	31.82052	63.65674	636.6192
2	0.288675	0.816497	1.885618	2.919986	4.30265	6.96456	9.92484	31.5991
3	0.276671	0.764892	1.637744	2.353363	3.18245	4.54070	5.84091	12.9240
4	0.270722	0.740697	1.533206	2.131847	2.77645	3.74695	4.60409	8.6103
5	0.267181	0.726687	1.475884	2.015048	2.57058	3.36493	4.03214	6.8688
6	0.264835	0.717558	1.439756	1.943180	2.44691	3.14267	3.70743	5.9588
7	0.263167	0.711142	1.414924	1.894579	2.36462	2.99795	3.49948	5.4079
8	0.261921	0.706387	1.396815	1.859548	2.30600	2.89646	3.35539	5.0413
9	0.260955	0.702722	1.383029	1.833113	2.26216	2.82144	3.24984	4.7809
10	0.260185	0.699812	1.372184	1.812461	2.22814	2.76377	3.16927	4.5869
11	0.259556	0.697445	1.363430	1.795885	2.20099	2.71808	3.10581	4.4370
12	0.259033	0.695483	1.356217	1.782288	2.17881	2.68100	3.05454	4.3178
13	0.258591	0.693829	1.350171	1.770933	2.16037	2.65031	3.01228	4.2208
14	0.258213	0.692417	1.345030	1.761310	2.14479	2.62449	2.97684	4.1405
15	0.257885	0.691197	1.340606	1.753050	2.13145	2.60248	2.94671	4.0728
16	0.257599	0.690132	1.336757	1.745884	2.11991	2.58349	2.92078	4.0150
17	0.257347	0.689195	1.333379	1.739607	2.10982	2.56693	2.89823	3.9651
18	0.257123	0.688364	1.330391	1.734064	2.10092	2.55238	2.87844	3.9216
19	0.256923	0.687621	1.327728	1.729133	2.09302	2.53948	2.86093	3.8834
20	0.256743	0.686954	1.325341	1.724718	2.08596	2.52798	2.84534	3.8495
21	0.256580	0.686352	1.323188	1.720743	2.07961	2.51765	2.83136	3.8193
22	0.256432	0.685805	1.321237	1.717144	2.07387	2.50832	2.81876	3.7921
23	0.256297	0.685306	1.319460	1.713872	2.06866	2.49987	2.80734	3.7676
24	0.256173	0.684850	1.317836	1.710882	2.06390	2.49216	2.79694	3.7454
25	0.256060	0.684430	1.316345	1.708141	2.05954	2.48511	2.78744	3.7251
26	0.255955	0.684043	1.314972	1.705618	2.05553	2.47863	2.77871	3.7066
27	0.255858	0.683685	1.313703	1.703288	2.05183	2.47266	2.77068	3.6896
28	0.255768	0.683353	1.312527	1.701131	2.04841	2.46714	2.76326	3.6739
29	0.255684	0.683044	1.311434	1.699127	2.04523	2.46202	2.75639	3.6594
30	0.255605	0.682756	1.310415	1.697261	2.04227	2.45726	2.75000	3.6460
∞	0.253347	0.674490	1.281552	1.644854	1.95996	2.32635	2.57583	3.2905

Binomial Table

Table A-2 shows probabilities for the binomial distribution (refer to Chapter 17). To use Table A-2, you need three pieces of information from the particular problem you're working on:

- ✓ The sample size, n
- ✓ The probability of success, p
- ✓ The value of x for which you want the cumulative probability

Find the portion of Table A-2 that's devoted to your n , and look at the row for your x and

the column for your p . Intersect that row and column, and you can see the probability for x . To get the probability of being strictly less than, greater than, greater than or equal to, or between two values of x , you sum the appropriate values of Table A-2, using the steps found in Chapter 16.

Table A-2

The Binomial Table

Numbers in the table represent the probabilities for values of x from 0 to n .

		p													
		0.1	0.2	0.25	0.3	0.4	0.5	0.6	0.7	0.75	0.8	0.9			
		$\binom{n}{x} p^x (1-p)^{n-x}$	n	x											
			1	0	0.900	0.800	0.750	0.700	0.600	0.500	0.400	0.300	0.250	0.200	0.100
			1	1	0.100	0.200	0.250	0.300	0.400	0.500	0.600	0.700	0.750	0.800	0.900
			2	0	0.810	0.640	0.563	0.490	0.360	0.250	0.160	0.090	0.063	0.040	0.010
			2	1	0.180	0.320	0.375	0.420	0.480	0.500	0.480	0.420	0.375	0.320	0.180
			2	2	0.010	0.040	0.063	0.090	0.160	0.250	0.360	0.490	0.563	0.640	0.810
			3	0	0.729	0.512	0.422	0.343	0.216	0.125	0.064	0.027	0.016	0.008	0.001
			3	1	0.243	0.384	0.422	0.441	0.432	0.375	0.288	0.189	0.141	0.096	0.027
			3	2	0.027	0.096	0.141	0.189	0.288	0.375	0.432	0.441	0.422	0.384	0.243
			3	3	0.001	0.008	0.016	0.027	0.064	0.125	0.216	0.343	0.422	0.512	0.729
			4	0	0.656	0.410	0.316	0.240	0.130	0.063	0.026	0.008	0.004	0.002	0.000
			4	1	0.292	0.410	0.422	0.412	0.346	0.250	0.154	0.076	0.047	0.026	0.004
			4	2	0.049	0.154	0.211	0.265	0.346	0.375	0.346	0.265	0.211	0.154	0.049
			4	3	0.004	0.026	0.047	0.076	0.154	0.250	0.346	0.412	0.422	0.410	0.292
			4	4	0.000	0.002	0.004	0.008	0.026	0.063	0.130	0.240	0.316	0.410	0.656
			5	0	0.590	0.328	0.237	0.168	0.078	0.031	0.010	0.002	0.001	0.000	0.000
			5	1	0.328	0.410	0.396	0.360	0.259	0.156	0.077	0.028	0.015	0.006	0.000
			5	2	0.073	0.205	0.264	0.309	0.346	0.312	0.230	0.132	0.088	0.051	0.008
			5	3	0.008	0.051	0.088	0.132	0.230	0.312	0.346	0.309	0.264	0.205	0.073
			5	4	0.000	0.006	0.015	0.028	0.077	0.156	0.259	0.360	0.396	0.410	0.328
			5	5	0.000	0.000	0.001	0.002	0.010	0.031	0.078	0.168	0.237	0.328	0.590
			6	0	0.531	0.262	0.178	0.118	0.047	0.016	0.004	0.001	0.000	0.000	0.000
			6	1	0.354	0.393	0.356	0.303	0.187	0.094	0.037	0.010	0.004	0.002	0.000
			6	2	0.098	0.246	0.297	0.324	0.311	0.234	0.138	0.060	0.033	0.015	0.001
			6	3	0.015	0.082	0.132	0.185	0.276	0.313	0.276	0.185	0.132	0.082	0.015
			6	4	0.001	0.015	0.033	0.060	0.138	0.234	0.311	0.324	0.297	0.246	0.098
			6	5	0.000	0.002	0.004	0.010	0.037	0.094	0.187	0.303	0.356	0.393	0.354
			6	6	0.000	0.000	0.000	0.001	0.004	0.016	0.047	0.118	0.178	0.262	0.531
			7	0	0.478	0.210	0.133	0.082	0.028	0.008	0.002	0.000	0.000	0.000	0.000
			7	1	0.372	0.367	0.311	0.247	0.131	0.055	0.017	0.004	0.001	0.000	0.000
			7	2	0.124	0.275	0.311	0.318	0.261	0.164	0.077	0.025	0.012	0.004	0.000
			7	3	0.023	0.115	0.173	0.227	0.290	0.273	0.194	0.097	0.058	0.029	0.003
			7	4	0.003	0.029	0.058	0.097	0.194	0.273	0.290	0.227	0.173	0.115	0.023
			7	5	0.000	0.004	0.012	0.025	0.077	0.164	0.261	0.318	0.311	0.275	0.124
			7	6	0.000	0.000	0.001	0.004	0.017	0.055	0.131	0.247	0.311	0.367	0.372
			7	7	0.000	0.000	0.000	0.000	0.002	0.008	0.028	0.082	0.133	0.210	0.478

Table A-2 (continued)

		p										
Binomial probabilities:		$\binom{n}{x} p^x (1-p)^{n-x}$										
n	x	0.1	0.2	0.25	0.3	0.4	0.5	0.6	0.7	0.75	0.8	0.9
8	0	0.430	0.168	0.100	0.058	0.017	0.004	0.001	0.000	0.000	0.000	0.000
	1	0.383	0.336	0.267	0.198	0.090	0.031	0.008	0.001	0.000	0.000	0.000
	2	0.149	0.294	0.311	0.296	0.209	0.109	0.041	0.010	0.004	0.001	0.000
	3	0.033	0.147	0.208	0.254	0.279	0.219	0.124	0.047	0.023	0.009	0.000
	4	0.005	0.046	0.087	0.136	0.232	0.273	0.232	0.136	0.087	0.046	0.005
	5	0.000	0.009	0.023	0.047	0.124	0.219	0.279	0.254	0.208	0.147	0.033
	6	0.000	0.001	0.004	0.010	0.041	0.109	0.209	0.296	0.311	0.294	0.149
	7	0.000	0.000	0.000	0.001	0.008	0.031	0.090	0.198	0.267	0.336	0.383
	8	0.000	0.000	0.000	0.000	0.001	0.004	0.017	0.058	0.100	0.168	0.430
9	0	0.387	0.134	0.075	0.040	0.010	0.002	0.000	0.000	0.000	0.000	0.000
	1	0.387	0.302	0.225	0.156	0.060	0.018	0.004	0.000	0.000	0.000	0.000
	2	0.172	0.302	0.300	0.267	0.161	0.070	0.021	0.004	0.001	0.000	0.000
	3	0.045	0.176	0.234	0.267	0.251	0.164	0.074	0.021	0.009	0.003	0.000
	4	0.007	0.066	0.117	0.172	0.251	0.246	0.167	0.074	0.039	0.017	0.001
	5	0.001	0.017	0.039	0.074	0.167	0.246	0.251	0.172	0.117	0.066	0.007
	6	0.000	0.003	0.009	0.021	0.074	0.164	0.251	0.267	0.234	0.176	0.045
	7	0.000	0.000	0.001	0.004	0.021	0.070	0.161	0.267	0.300	0.302	0.172
	8	0.000	0.000	0.000	0.000	0.004	0.018	0.060	0.156	0.225	0.302	0.387
10	0	0.349	0.107	0.056	0.028	0.006	0.001	0.000	0.000	0.000	0.000	0.000
	1	0.387	0.268	0.188	0.121	0.040	0.010	0.002	0.000	0.000	0.000	0.000
	2	0.194	0.302	0.282	0.233	0.121	0.044	0.011	0.001	0.000	0.000	0.000
	3	0.057	0.201	0.250	0.267	0.215	0.117	0.042	0.009	0.003	0.001	0.000
	4	0.011	0.088	0.146	0.200	0.251	0.205	0.111	0.037	0.016	0.006	0.000
	5	0.001	0.026	0.058	0.103	0.201	0.246	0.201	0.103	0.058	0.026	0.001
	6	0.000	0.006	0.016	0.037	0.111	0.205	0.251	0.200	0.146	0.088	0.011
	7	0.000	0.001	0.003	0.009	0.042	0.117	0.215	0.267	0.250	0.201	0.057
	8	0.000	0.000	0.000	0.001	0.011	0.044	0.121	0.233	0.282	0.302	0.194
11	0	0.300	0.086	0.042	0.020	0.004	0.000	0.000	0.000	0.000	0.000	0.000
	1	0.384	0.236	0.155	0.093	0.027	0.005	0.001	0.000	0.000	0.000	0.000
	2	0.213	0.295	0.258	0.200	0.089	0.027	0.005	0.001	0.000	0.000	0.000
	3	0.071	0.221	0.258	0.257	0.177	0.081	0.023	0.004	0.001	0.000	0.000
	4	0.016	0.111	0.172	0.220	0.236	0.161	0.070	0.017	0.006	0.002	0.000
	5	0.002	0.039	0.080	0.132	0.221	0.226	0.147	0.057	0.027	0.010	0.000
	6	0.000	0.010	0.027	0.057	0.147	0.226	0.221	0.132	0.080	0.039	0.002
	7	0.000	0.002	0.006	0.017	0.070	0.161	0.236	0.220	0.172	0.111	0.016
	8	0.000	0.000	0.001	0.004	0.023	0.081	0.177	0.257	0.258	0.221	0.071
12	0	0.289	0.086	0.042	0.020	0.004	0.000	0.000	0.000	0.000	0.000	0.000
	1	0.380	0.236	0.155	0.093	0.027	0.005	0.001	0.000	0.000	0.000	0.000
	2	0.212	0.295	0.258	0.200	0.089	0.027	0.005	0.001	0.000	0.000	0.000

Table A-2 (continued)

		p												
Binomial probabilities:		n	x	0.1	0.2	0.25	0.3	0.4	0.5	0.6	0.7	0.75	0.8	0.9
$\binom{n}{x}$	$p^x(1-p)^{n-x}$	12	0	0.282	0.069	0.032	0.014	0.002	0.000	0.000	0.000	0.000	0.000	0.000
		1	1	0.377	0.206	0.127	0.071	0.017	0.003	0.000	0.000	0.000	0.000	0.000
		2	2	0.230	0.283	0.232	0.168	0.064	0.016	0.002	0.000	0.000	0.000	0.000
		3	3	0.085	0.236	0.258	0.240	0.142	0.054	0.012	0.001	0.000	0.000	0.000
		4	4	0.021	0.133	0.194	0.231	0.213	0.121	0.042	0.008	0.002	0.001	0.000
		5	5	0.004	0.053	0.103	0.158	0.227	0.193	0.101	0.029	0.011	0.003	0.000
		6	6	0.000	0.016	0.040	0.079	0.177	0.226	0.177	0.079	0.040	0.016	0.000
		7	7	0.000	0.003	0.011	0.029	0.101	0.193	0.227	0.158	0.103	0.053	0.004
		8	8	0.000	0.001	0.002	0.008	0.042	0.121	0.213	0.231	0.194	0.133	0.021
		9	9	0.000	0.000	0.000	0.001	0.012	0.054	0.142	0.240	0.258	0.236	0.085
		10	10	0.000	0.000	0.000	0.000	0.002	0.016	0.064	0.168	0.232	0.283	0.230
		11	11	0.000	0.000	0.000	0.000	0.000	0.003	0.017	0.071	0.127	0.206	0.377
		12	12	0.000	0.000	0.000	0.000	0.000	0.000	0.002	0.014	0.032	0.069	0.282
13	0	13	0	0.254	0.055	0.024	0.010	0.001	0.000	0.000	0.000	0.000	0.000	0.000
	1	1	1	0.367	0.179	0.103	0.054	0.011	0.002	0.000	0.000	0.000	0.000	0.000
	2	2	2	0.245	0.268	0.206	0.139	0.045	0.010	0.001	0.000	0.000	0.000	0.000
	3	3	3	0.100	0.246	0.252	0.218	0.111	0.035	0.006	0.001	0.000	0.000	0.000
	4	4	4	0.028	0.154	0.210	0.234	0.184	0.087	0.024	0.003	0.001	0.000	0.000
	5	5	5	0.006	0.069	0.126	0.180	0.221	0.157	0.066	0.014	0.005	0.001	0.000
	6	6	6	0.001	0.023	0.056	0.103	0.197	0.209	0.131	0.044	0.019	0.006	0.000
	7	7	7	0.000	0.006	0.019	0.044	0.131	0.209	0.197	0.103	0.056	0.023	0.001
	8	8	8	0.000	0.001	0.005	0.014	0.066	0.157	0.221	0.180	0.126	0.069	0.006
	9	9	9	0.000	0.000	0.001	0.003	0.024	0.087	0.184	0.234	0.210	0.154	0.028
	10	10	10	0.000	0.000	0.000	0.001	0.006	0.035	0.111	0.218	0.252	0.246	0.100
	11	11	11	0.000	0.000	0.000	0.000	0.001	0.010	0.045	0.139	0.206	0.268	0.245
	12	12	12	0.000	0.000	0.000	0.000	0.002	0.011	0.054	0.103	0.179	0.367	
	13	13	13	0.000	0.000	0.000	0.000	0.000	0.001	0.010	0.024	0.055	0.254	
14	0	14	0	0.229	0.044	0.018	0.007	0.001	0.000	0.000	0.000	0.000	0.000	0.000
	1	1	1	0.355	0.154	0.083	0.041	0.007	0.001	0.000	0.000	0.000	0.000	0.000
	2	2	2	0.257	0.250	0.180	0.113	0.032	0.006	0.001	0.000	0.000	0.000	0.000
	3	3	3	0.114	0.250	0.240	0.194	0.085	0.022	0.003	0.000	0.000	0.000	0.000
	4	4	4	0.035	0.172	0.220	0.229	0.155	0.061	0.014	0.001	0.000	0.000	0.000
	5	5	5	0.008	0.086	0.147	0.196	0.207	0.122	0.041	0.007	0.002	0.000	0.000
	6	6	6	0.001	0.032	0.073	0.126	0.207	0.183	0.092	0.023	0.008	0.002	0.000
	7	7	7	0.000	0.009	0.028	0.062	0.157	0.209	0.157	0.062	0.028	0.009	0.000
	8	8	8	0.000	0.002	0.008	0.023	0.092	0.183	0.207	0.126	0.073	0.032	0.001
	9	9	9	0.000	0.000	0.002	0.007	0.041	0.122	0.207	0.196	0.147	0.086	0.008
	10	10	10	0.000	0.000	0.000	0.001	0.014	0.061	0.155	0.229	0.220	0.172	0.035
	11	11	11	0.000	0.000	0.000	0.000	0.003	0.022	0.085	0.194	0.240	0.250	0.114
	12	12	12	0.000	0.000	0.000	0.000	0.001	0.006	0.032	0.113	0.180	0.250	0.257
	13	13	13	0.000	0.000	0.000	0.000	0.000	0.001	0.007	0.041	0.083	0.154	0.356
	14	14	14	0.000	0.000	0.000	0.000	0.000	0.001	0.007	0.018	0.044	0.229	

Table A-2 (continued)

		p											
		0.1	0.2	0.25	0.3	0.4	0.5	0.6	0.7	0.75	0.8	0.9	
Binomial probabilities:		$\binom{n}{x} p^x (1-p)^{n-x}$	n	x									
$\binom{n}{x} p^x (1-p)^{n-x}$	15	0	0.206	0.035	0.013	0.005	0.000	0.000	0.000	0.000	0.000	0.000	
		1	0.343	0.132	0.067	0.031	0.005	0.000	0.000	0.000	0.000	0.000	
		2	0.267	0.231	0.156	0.092	0.022	0.003	0.000	0.000	0.000	0.000	
		3	0.129	0.250	0.225	0.170	0.063	0.014	0.002	0.000	0.000	0.000	
		4	0.043	0.188	0.225	0.219	0.127	0.042	0.007	0.001	0.000	0.000	
		5	0.010	0.103	0.165	0.206	0.186	0.092	0.024	0.003	0.001	0.000	
		6	0.002	0.043	0.092	0.147	0.207	0.153	0.061	0.012	0.003	0.001	0.000
		7	0.000	0.014	0.039	0.081	0.177	0.196	0.118	0.035	0.013	0.003	0.000
		8	0.000	0.003	0.013	0.035	0.118	0.196	0.177	0.081	0.039	0.014	0.000
		9	0.000	0.001	0.003	0.012	0.061	0.153	0.207	0.147	0.092	0.043	0.002
		10	0.000	0.000	0.001	0.003	0.024	0.092	0.186	0.206	0.165	0.103	0.010
		11	0.000	0.000	0.000	0.001	0.007	0.042	0.127	0.219	0.225	0.188	0.043
		12	0.000	0.000	0.000	0.000	0.002	0.014	0.063	0.170	0.225	0.250	0.129
		13	0.000	0.000	0.000	0.000	0.000	0.003	0.022	0.092	0.156	0.231	0.267
		14	0.000	0.000	0.000	0.000	0.000	0.000	0.005	0.031	0.067	0.132	0.343
		15	0.000	0.000	0.000	0.000	0.000	0.000	0.005	0.013	0.035	0.206	
$\binom{n}{x} p^x (1-p)^{n-x}$	20	0	0.122	0.012	0.003	0.001	0.000	0.000	0.000	0.000	0.000	0.000	
		1	0.270	0.058	0.021	0.007	0.000	0.000	0.000	0.000	0.000	0.000	
		2	0.285	0.137	0.067	0.028	0.003	0.000	0.000	0.000	0.000	0.000	
		3	0.190	0.205	0.134	0.072	0.012	0.001	0.000	0.000	0.000	0.000	
		4	0.090	0.218	0.190	0.130	0.035	0.005	0.000	0.000	0.000	0.000	
		5	0.032	0.175	0.202	0.179	0.075	0.015	0.001	0.000	0.000	0.000	
		6	0.009	0.109	0.169	0.192	0.124	0.037	0.005	0.000	0.000	0.000	
		7	0.002	0.055	0.112	0.164	0.166	0.074	0.015	0.001	0.000	0.000	
		8	0.000	0.022	0.061	0.114	0.180	0.120	0.035	0.004	0.001	0.000	
		9	0.000	0.007	0.027	0.065	0.160	0.160	0.071	0.012	0.003	0.000	
		10	0.000	0.002	0.010	0.031	0.117	0.176	0.117	0.031	0.010	0.002	
		11	0.000	0.000	0.003	0.012	0.071	0.160	0.160	0.065	0.027	0.007	
		12	0.000	0.000	0.001	0.004	0.035	0.120	0.180	0.114	0.061	0.022	
		13	0.000	0.000	0.000	0.001	0.015	0.074	0.166	0.164	0.112	0.055	
		14	0.000	0.000	0.000	0.000	0.005	0.037	0.124	0.192	0.169	0.109	
		15	0.000	0.000	0.000	0.000	0.001	0.015	0.075	0.179	0.202	0.175	
		16	0.000	0.000	0.000	0.000	0.000	0.005	0.035	0.130	0.190	0.218	
		17	0.000	0.000	0.000	0.000	0.000	0.001	0.012	0.072	0.134	0.205	
		18	0.000	0.000	0.000	0.000	0.000	0.003	0.028	0.067	0.137	0.285	
		19	0.000	0.000	0.000	0.000	0.000	0.000	0.007	0.021	0.058	0.270	
		20	0.000	0.000	0.000	0.000	0.000	0.001	0.003	0.012	0.122		

Chi-Square Table

Table A-3 shows right-tail probabilities for the Chi-square distribution (you can use Chapter 14 as a reference for the Chi-square test). To use Table A-3, you need three pieces of information from the particular problem you're working on:

- ✓ The sample size, n .
- ✓ The value of Chi-squared for which you want the right-tail probability.
- ✓ If you're working with a two-way table, you need r = number of rows and c = number of columns. If you're working with a goodness-of-fit test, you need $k - 1$, where k is the number of categories.

The degrees of freedom for the Chi-square test statistic is $(r - 1) * (c - 1)$ if you're testing for an association between two variables, where r and c are the number of rows and columns in the two-way table, respectively. Or, the degrees of freedom is $k - 1$ in a goodness-of-fit test, where k is the number of categories; see Chapter 15.

Go across the row for your degrees of freedom until you find the value in that row closest to your Chi-square test statistic. Look up at the number at the top of that column.

That value is the area to the right of (beyond) that particular Chi-square statistic.

Table A-3

The Chi-Square Table

Numbers in the table represent Chi-square values whose area to the right equals p .

$\frac{df}{p}$	0.10	0.05	0.025	0.01	0.005
1	2.71	3.84	5.02	6.64	7.88
2	4.61	5.99	7.38	9.21	10.60
3	6.25	7.82	9.35	11.35	12.84
4	7.78	9.49	11.14	13.28	14.86
5	9.24	11.07	12.83	15.09	16.75
6	10.65	12.59	14.45	16.81	18.55
7	12.02	14.07	16.01	18.48	20.28
8	13.36	15.51	17.54	20.09	21.96
9	14.68	16.92	19.02	21.67	23.59
10	15.99	18.31	20.48	23.21	25.19
11	17.28	19.68	21.92	24.73	26.76
12	18.55	21.03	23.34	26.22	28.30
13	19.81	22.36	24.74	27.69	29.819
14	21.06	23.69	26.12	29.14	31.32
15	22.31	25.00	27.49	30.58	32.80
16	23.54	26.30	28.85	32.00	34.27
17	24.77	27.59	30.19	33.41	35.72
18	25.99	28.87	31.53	34.81	37.16
19	27.20	30.14	32.85	36.19	38.58
20	28.41	31.41	34.17	37.57	40.00
21	29.62	32.67	35.48	38.93	41.40
22	30.81	33.92	36.78	40.29	42.80
23	32.01	35.17	38.08	41.64	44.18
24	33.20	36.42	39.36	42.98	45.56
25	34.38	37.65	40.65	44.31	46.93
26	35.56	38.89	41.92	45.64	48.29
27	36.74	40.11	43.20	46.96	49.65
28	37.92	41.34	44.46	48.28	50.99
29	39.09	42.56	45.72	49.59	52.34
30	40.26	43.77	46.98	50.89	53.67
40	51.81	55.76	59.34	63.69	66.77
50	63.17	67.51	71.42	76.15	79.49

Rank Sum Table

Tables A-4(a) and A-4(b) show the critical values for the rank sum test for $\alpha = 0.05$ or $\alpha = 0.10$, $\rho\varepsilon\sigma\varepsilon\chi\tau\omega\lambda\psi$; see Chapter 18 for more on this test. To use Table A-4, you need two pieces of information from the particular problem you're working on:

- ✓ The rank sum statistic, T
- ✓ The sample sizes of the two samples, n_1 and n_2

To find the critical value for your rank sum statistic using Table A-4, go to the column

representing n_1 and the row representing n_2 . Intersect the row and column to find the lower and upper critical values (denoted T_L and T_U) for the rank sum test.

Table A-4(a)

The Rank Sum Table ($\alpha = 0.05$)

$n_1 \backslash n_2$	3	4	5	6	7	8	9	10
	T_L	T_U	T_L	T_U	T_L	T_U	T_L	T_U
3	5	16	6	18	6	21	7	23
4	6	18	11	25	12	28	13	35
5	6	21	12	28	18	37	19	41
6	7	23	12	32	19	41	26	52
7	7	26	13	35	20	45	28	56
8	8	28	14	38	21	49	29	61
9	8	31	15	41	22	53	31	65
10	9	33	16	44	24	56	32	70

Table A-4(b)

The Rank Sum Table ($\alpha = 0.10$)

$n_1 \backslash n_2$	3	4	5	6	7	8	9	10
	T_L	T_U	T_L	T_U	T_L	T_U	T_L	T_U
3	6	15	7	17	7	20	8	22
4	7	17	12	24	13	27	14	30
5	7	20	13	37	19	36	20	40
6	8	22	14	30	20	40	28	50
7	9	24	15	33	22	43	30	54
8	9	27	16	36	24	46	32	58
9	10	29	17	39	25	50	33	63
10	11	31	18	42	26	54	35	67

F-Table

Table A-5 shows the critical values on the F -distribution where α is equal to 0.05. (*Critical values* are those values that represent the boundary between rejecting H_0 and not rejecting H_0 ; refer to Chapter 9.) To use Table A-5, you need three pieces of information from the particular problem you're working on:

- ✓ The sample size, n
- ✓ The number of populations (or treatments being compared), k
- ✓ The value of F for which you want the cumulative probability

To find the critical value for your F -test statistic using Table A-5, go to the column representing the degrees of freedom you need ($k - 1$ and $n - k$). Intersect the column degrees of freedom ($k - 1$) with the row degrees of freedom ($n - k$), and you find the critical value on the F -distribution.

Table A-5

The F-Table ($\alpha = 0.05$)

		F [(df_1, df_2)]																	
		1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	40	60	120
df2/df1																			
1	161.4476	199.5000	215.7073	224.5832	230.1619	233.9860	236.7684	238.8827	240.5433	241.8817	243.9060	245.9499	248.0131	249.0518	250.0951	251.1432	252.1957	253.252	
2	18.5128	19.0000	19.1643	19.2468	19.2964	19.3295	19.3532	19.3710	19.3848	19.4021	19.4291	19.4458	19.4651	19.4824	19.4970	19.4971	19.4987		
3	10.1260	9.5521	9.2706	9.1172	9.0135	8.9406	8.8867	8.8450	8.8123	8.7855	8.7446	8.7029	8.6652	8.6385	8.6166	8.5944	8.5720	8.549	
4	7.7098	6.9443	6.5914	6.3882	6.1951	6.0942	6.0410	5.9983	5.9444	5.8944	5.8451	5.7944	5.7444	5.7449	5.7170	5.6877	5.658		
5	5.6079	5.7861	5.4095	5.1922	5.0503	4.9503	4.8759	4.8183	4.7725	4.7351	4.6777	4.6188	4.5581	4.5272	4.4957	4.4638	4.4314	4.388	
6	5.3674	5.1433	4.7571	4.5337	4.3074	4.2639	4.2067	4.1468	4.0990	4.0500	3.9999	3.9381	3.8742	3.8415	3.8082	3.7743	3.7388	3.704	
7	5.5914	4.7374	4.2468	4.1203	3.9715	3.8650	3.7870	3.7257	3.6767	3.6265	3.5747	3.5107	3.4445	3.4105	3.3758	3.3404	3.3043	3.267	
8	5.3177	4.4590	4.0662	3.8379	3.6875	3.5806	3.5005	3.4381	3.3881	3.3289	3.2839	3.2384	3.1953	3.1152	3.0794	3.0428	3.0053	2.966	
9	5.1174	4.2565	3.6625	3.6331	3.4817	3.3738	3.2927	3.2296	3.1789	3.1373	3.0729	3.0061	2.9365	2.8657	2.8055	2.7459	2.7147		
10	4.9546	4.1028	3.7083	3.4780	3.3258	3.2172	3.1555	3.0717	3.0204	2.9782	2.9130	2.8450	2.7740	2.7372	2.6996	2.6609	2.6211	2.580	
11	4.8443	3.9823	3.5604	3.3567	3.2039	3.0946	3.0123	2.9480	2.8962	2.8536	2.8186	2.6454	2.6090	2.5705	2.5309	2.4901	2.448		
12	4.7472	3.8853	3.4903	3.2992	3.1059	2.9961	2.9134	2.8486	2.7964	2.7534	2.6866	2.6169	2.5436	2.5055	2.4683	2.4259	2.3842	2.341	
13	4.6672	3.8056	3.4105	3.1791	3.0254	2.9153	2.8231	2.7669	2.7144	2.6710	2.6037	2.5331	2.4589	2.4202	2.3893	2.3392	2.2966	2.252	
14	4.6001	3.7389	3.3439	3.1122	2.9582	2.8487	2.7642	2.6887	2.6458	2.6022	2.5342	2.4652	2.3879	2.3487	2.3062	2.2664	2.2229	2.177	
15	4.5431	3.6823	3.2874	3.0556	2.9013	2.7995	2.7066	2.6468	2.5876	2.5437	2.4753	2.4034	2.3278	2.2878	2.2468	2.2043	2.1601	2.114	
16	4.4940	3.6337	3.2389	3.0069	2.8524	2.7413	2.6572	2.5911	2.5377	2.4935	2.4247	2.3522	2.2756	2.2254	2.1938	2.1507	2.1058	2.058	
17	4.4513	3.5915	3.1968	2.9647	2.8100	2.6887	2.6143	2.5480	2.4943	2.4499	2.3807	2.3077	2.2304	2.1886	2.1477	2.1040	2.0584	2.010	
18	4.4139	3.5546	3.1599	2.9277	2.7729	2.6613	2.5767	2.5102	2.4563	2.4117	2.3421	2.2886	2.1906	2.1487	2.1071	2.0629	2.0166	1.968	
19	4.3807	3.5219	3.1724	2.9551	2.7401	2.6283	2.5455	2.4768	2.4227	2.3779	2.3080	2.2341	2.1555	2.1141	2.0712	2.0264	1.9795	1.930	
20	4.3512	3.4926	3.0984	2.8561	2.7109	2.5980	2.5140	2.4471	2.3938	2.3479	2.2776	2.2033	2.1242	2.065	2.035	1.9838	1.9454	1.8956	
21	4.3248	3.4668	3.0725	2.8401	2.6848	2.5727	2.4805	2.4265	2.3660	2.3210	2.2504	2.1757	2.0960	2.0540	2.0102	1.9645	1.9165	1.885	
22	4.3009	3.4434	3.0491	2.8167	2.6613	2.5491	2.4638	2.3965	2.3419	2.2967	2.2258	2.1508	2.0707	2.0263	1.9842	1.9380	1.8894	1.838	
23	4.2793	3.4221	3.0280	2.7955	2.6400	2.5277	2.4422	2.3748	2.3201	2.2747	2.2036	2.1282	2.0476	2.0050	1.9605	1.9139	1.8648	1.812	
24	4.2597	3.4028	3.0088	2.7763	2.6283	2.5140	2.4471	2.3938	2.3479	2.2776	2.2033	2.1242	2.065	2.0267	1.9838	1.9390	1.8920	1.8424	
25	4.2417	3.3852	2.9912	2.7587	2.6030	2.4904	2.4247	2.3371	2.2821	2.2365	2.1649	2.0889	2.0375	1.9643	1.9192	1.8718	1.8217	1.768	
26	4.2252	3.3690	2.9752	2.7426	2.5868	2.4741	2.3883	2.3265	2.2655	2.2197	2.1479	2.0716	1.9898	1.9464	1.9010	1.8533	1.8027	1.7448	
27	4.2100	3.3541	2.9604	2.7278	2.5719	2.4951	2.3752	2.3053	2.2501	2.2043	2.1323	2.0558	1.9736	1.9299	1.8842	1.8361	1.7851	1.730	
28	4.1980	3.3404	2.9467	2.7141	2.5581	2.4453	2.3693	2.2913	2.2360	2.1900	2.1179	2.0411	1.9586	1.9147	1.8687	1.8203	1.7689	1.713	
29	4.1830	3.3277	2.9340	2.7014	2.5454	2.4324	2.3463	2.2783	2.2229	2.1768	2.1045	2.0275	1.9446	1.9005	1.8543	1.8055	1.7537	1.698	
30	4.1709	3.3158	2.9223	2.6896	2.5338	2.4265	2.3343	2.2662	2.2107	2.1646	2.0921	2.0148	1.9317	1.8874	1.8409	1.7918	1.7396	1.683	
40	4.0847	3.2317	2.8387	2.6060	2.4495	2.3359	2.2490	2.1802	2.1240	2.0772	2.0035	1.9245	1.8389	1.7929	1.7444	1.6928	1.6373	1.576	
60	4.0012	3.1504	2.7581	2.5552	2.3683	2.2541	2.1665	2.0970	2.0401	1.9526	1.9174	1.8484	1.7480	1.7001	1.6491	1.5943	1.5343	1.467	
120	3.9201	3.0718	2.6802	2.4472	2.2899	2.1750	2.0688	2.0164	1.9688	1.9105	1.8337	1.7955	1.6587	1.6084	1.5543	1.4952	1.4290	1.351	



Get More and Do More at Dummies.com®



Start with **FREE** Cheat Sheets

Cheat Sheets include

- Checklists
- Charts
- Common Instructions
- And Other Good Stuff!

To access the cheat sheet specifically for this book, go to
www.dummies.com/cheatsheet/statistics2.

Get Smart at Dummies.com

Dummies.com makes your life easier with 1,000s of answers on everything from removing wallpaper to using the latest version of Windows.

Check out our

- Videos
- Illustrated Articles
- Step-by-Step Instructions

Plus, each month you can win valuable prizes by entering our Dummies.com sweepstakes.*

Want a weekly dose of Dummies? Sign up for Newsletters on

- Digital Photography
- Microsoft Windows & Office
- Personal Finance & Investing
- Health & Wellness
- Computing, iPods & Cell Phones
- eBay
- Internet
- Food, Home & Garden



*Sweepstakes not currently available in all countries; visit Dummies.com for official rules.

Find out "HOW" at Dummies.com