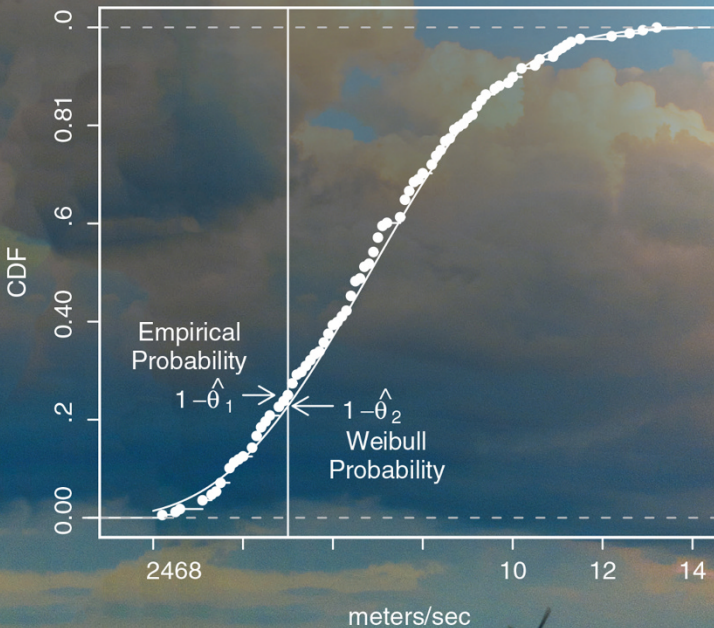


SECOND EDITION

MATHEMATICAL STATISTICS WITH RESAMPLING AND R

LAURA M. CHIHARA | TIM C. HESTERBERG



WILEY

Mathematical Statistics with Resampling and R

Mathematical Statistics with Resampling and R

Second Edition

Laura M. Chihara
Carleton College

Tim C. Hesterberg
Google

WILEY

This second edition first published 2019
© 2019 by John Wiley & Sons, Inc.

Edition History

Mathematical Statistics with Resampling and R, Wiley, 2011

All rights reserved. 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 or otherwise, except as permitted by law. Advice on how to obtain permission to reuse material from this title is available at <http://www.wiley.com/go/permissions>.

The right of Laura M. Chihara and Tim C. Hesterberg to be identified as the author(s) of this work has been asserted in accordance with law.

Registered Office

John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, USA

Editorial Office

111 River Street, Hoboken, NJ 07030, USA

For details of our global editorial offices, customer services, and more information about Wiley products visit us at www.wiley.com.

Wiley also publishes its books in a variety of electronic formats and by print-on-demand. Some content that appears in standard print versions of this book may not be available in other formats.

Limit of Liability/Disclaimer of Warranty

The publisher and the authors 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 any implied warranties of fitness for a particular purpose. This work is sold with the understanding that the publisher is not engaged in rendering professional services. The advice and strategies contained herein may not be suitable for every situation. In view of on-going research, equipment modifications, changes in governmental regulations, and the constant flow of information relating to the use of experimental reagents, equipment, and devices, the reader is urged to review and evaluate the information provided in the package insert or instructions for each chemical, piece of equipment, reagent, or device for, among other things, any changes in the instructions or indication of usage and for added warnings and precautions. The fact that an organization or website is referred to in this work as a citation and/or 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 websites listed in this work may have changed or disappeared between when this work was written and when it is read. No warranty may be created or extended by any promotional statements for this work. Neither the publisher nor the author shall be liable for any damages arising here from.

Library of Congress Cataloging-in-Publication Data:

Names: Chihara, Laura, 1957– author. | Hesterberg, Tim, 1959– author.

Title: Mathematical statistics with resampling and R / Laura M. Chihara
(Carleton College), Tim C. Hesterberg (Google).

Description: Second edition. | Hoboken, NJ : Wiley, 2019. | Includes
bibliographical references and index. |

Identifiers: LCCN 2018008774 (print) | LCCN 2018013587 (ebook) | ISBN
9781119416524 (pdf) | ISBN 9781119416531 (epub) | ISBN 9781119416548
(cloth)

Subjects: LCSH: Resampling (Statistics) | Statistics. | Statistics–Data
processing. | Mathematical statistics–Data processing. | R (Computer
program language)

Classification: LCC QA278.8 (ebook) | LCC QA278.8 .C45 2018 (print) | DDC
519.5/4–dc23

LC record available at <https://lcn.loc.gov/2018008774>

Cover design by Wiley

Cover images: Courtesy of Carleton College

Set in 10/12pt Warnock by SPi Global, Pondicherry, India

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

*The world seldom notices who teachers are;
but civilization depends on what they do.*
– Lindley Stiles

To:
Theodore S. Chihara

To:
Bev Hesterberg

Contents

Preface *xiii*

1	Data and Case Studies	1
1.1	Case Study: Flight Delays	1
1.2	Case Study: Birth Weights of Babies	2
1.3	Case Study: Verizon Repair Times	3
1.4	Case Study: Iowa Recidivism	4
1.5	Sampling	5
1.6	Parameters and Statistics	6
1.7	Case Study: General Social Survey	7
1.8	Sample Surveys	8
1.9	Case Study: Beer and Hot Wings	9
1.10	Case Study: Black Spruce Seedlings	10
1.11	Studies	10
1.12	Google Interview Question: Mobile Ads Optimization	12
	Exercises	16
2	Exploratory Data Analysis	21
2.1	Basic Plots	21
2.2	Numeric Summaries	25
2.2.1	Center	25
2.2.2	Spread	26
2.2.3	Shape	27
2.3	Boxplots	28
2.4	Quantiles and Normal Quantile Plots	29
2.5	Empirical Cumulative Distribution Functions	35
2.6	Scatter Plots	38
2.7	Skewness and Kurtosis	40
	Exercises	42

3	Introduction to Hypothesis Testing: Permutation Tests	47
3.1	Introduction to Hypothesis Testing	47
3.2	Hypotheses	48
3.3	Permutation Tests	50
3.3.1	Implementation Issues	55
3.3.2	One-sided and Two-sided Tests	61
3.3.3	Other Statistics	62
3.3.4	Assumptions	64
3.3.5	Remark on Terminology	68
3.4	Matched Pairs	68
	Exercises	70
4	Sampling Distributions	75
4.1	Sampling Distributions	75
4.2	Calculating Sampling Distributions	80
4.3	The Central Limit Theorem	84
4.3.1	CLT for Binomial Data	86
4.3.2	Continuity Correction for Discrete Random Variables	89
4.3.3	Accuracy of the Central Limit Theorem*	91
4.3.4	CLT for Sampling Without Replacement	92
	Exercises	93
5	Introduction to Confidence Intervals: The Bootstrap	103
5.1	Introduction to the Bootstrap	103
5.2	The Plug-in Principle	110
5.2.1	Estimating the Population Distribution	112
5.2.2	How Useful Is the Bootstrap Distribution?	113
5.3	Bootstrap Percentile Intervals	118
5.4	Two-Sample Bootstrap	119
5.4.1	Matched Pairs	124
5.5	Other Statistics	128
5.6	Bias	131
5.7	Monte Carlo Sampling: The “Second Bootstrap Principle”	134
5.8	Accuracy of Bootstrap Distributions	135
5.8.1	Sample Mean: Large Sample Size	135
5.8.2	Sample Mean: Small Sample Size	137
5.8.3	Sample Median	138
5.8.4	Mean–Variance Relationship	138
5.9	How Many Bootstrap Samples Are Needed?	140
	Exercises	141
6	Estimation	149
6.1	Maximum Likelihood Estimation	149

6.1.1	Maximum Likelihood for Discrete Distributions	150
6.1.2	Maximum Likelihood for Continuous Distributions	153
6.1.3	Maximum Likelihood for Multiple Parameters	157
6.2	Method of Moments	161
6.3	Properties of Estimators	163
6.3.1	Unbiasedness	164
6.3.2	Efficiency	167
6.3.3	Mean Square Error	171
6.3.4	Consistency	173
6.3.5	Transformation Invariance*	175
6.3.6	Asymptotic Normality of MLE*	177
6.4	Statistical Practice	178
6.4.1	Are You Asking the Right Question?	179
6.4.2	Weights	179
	Exercises	180
7	More Confidence Intervals	187
7.1	Confidence Intervals for Means	187
7.1.1	Confidence Intervals for a Mean, Variance Known	187
7.1.2	Confidence Intervals for a Mean, Variance Unknown	192
7.1.3	Confidence Intervals for a Difference in Means	198
7.1.4	Matched Pairs, Revisited	204
7.2	Confidence Intervals in General	204
7.2.1	Location and Scale Parameters*	208
7.3	One-sided Confidence Intervals	212
7.4	Confidence Intervals for Proportions	214
7.4.1	Agresti–Coull Intervals for a Proportion	217
7.4.2	Confidence Intervals for a Difference of Proportions	218
7.5	Bootstrap Confidence Intervals	219
7.5.1	t Confidence Intervals Using Bootstrap Standard Errors	219
7.5.2	Bootstrap t Confidence Intervals	220
7.5.3	Comparing Bootstrap t and Formula t Confidence Intervals	224
7.6	Confidence Interval Properties	226
7.6.1	Confidence Interval Accuracy	226
7.6.2	Confidence Interval Length	227
7.6.3	Transformation Invariance	227
7.6.4	Ease of Use and Interpretation	227
7.6.5	Research Needed	228
	Exercises	228
8	More Hypothesis Testing	241
8.1	Hypothesis Tests for Means and Proportions: One Population	241
8.1.1	A Single Mean	241

8.1.2	One Proportion	244
8.2	Bootstrap t -Tests	246
8.3	Hypothesis Tests for Means and Proportions: Two Populations	248
8.3.1	Comparing Two Means	248
8.3.2	Comparing Two Proportions	251
8.3.3	Matched Pairs for Proportions	254
8.4	Type I and Type II Errors	255
8.4.1	Type I Errors	257
8.4.2	Type II Errors and Power	261
8.4.3	P -Values Versus Critical Regions	266
8.5	Interpreting Test Results	267
8.5.1	P -Values	267
8.5.2	On Significance	268
8.5.3	Adjustments for Multiple Testing	269
8.6	Likelihood Ratio Tests	271
8.6.1	Simple Hypotheses and the Neyman–Pearson Lemma	271
8.6.2	Likelihood Ratio Tests for Composite Hypotheses	275
8.7	Statistical Practice	279
8.7.1	More Campaigns with No Clicks and No Conversions	284
	Exercises	285
9	Regression	297
9.1	Covariance	297
9.2	Correlation	301
9.3	Least-Squares Regression	304
9.3.1	Regression Toward the Mean	308
9.3.2	Variation	310
9.3.3	Diagnostics	311
9.3.4	Multiple Regression	317
9.4	The Simple Linear Model	317
9.4.1	Inference for α and β	322
9.4.2	Inference for the Response	326
9.4.3	Comments About Assumptions for the Linear Model	330
9.5	Resampling Correlation and Regression	332
9.5.1	Permutation Tests	335
9.5.2	Bootstrap Case Study: Bushmeat	336
9.6	Logistic Regression	340
9.6.1	Inference for Logistic Regression	346
	Exercises	350
10	Categorical Data	359
10.1	Independence in Contingency Tables	359
10.2	Permutation Test of Independence	361

10.3	Chi-square Test of Independence	365
10.3.1	Model for Chi-square Test of Independence	366
10.3.2	2×2 Tables	368
10.3.3	Fisher's Exact Test	370
10.3.4	Conditioning	371
10.4	Chi-square Test of Homogeneity	372
10.5	Goodness-of-fit Tests	374
10.5.1	All Parameters Known	374
10.5.2	Some Parameters Estimated	377
10.6	Chi-square and the Likelihood Ratio*	379
	Exercises	380
11	Bayesian Methods	391
11.1	Bayes Theorem	392
11.2	Binomial Data: Discrete Prior Distributions	392
11.3	Binomial Data: Continuous Prior Distributions	400
11.4	Continuous Data	406
11.5	Sequential Data	409
	Exercises	414
12	One-way ANOVA	419
12.1	Comparing Three or More Populations	419
12.1.1	The ANOVA F-test	419
12.1.2	A Permutation Test Approach	428
	Exercises	429
13	Additional Topics	433
13.1	Smoothed Bootstrap	433
13.1.1	Kernel Density Estimate	435
13.2	Parametric Bootstrap	437
13.3	The Delta Method	441
13.4	Stratified Sampling	445
13.5	Computational Issues in Bayesian Analysis	446
13.6	Monte Carlo Integration	448
13.7	Importance Sampling	452
13.7.1	Ratio Estimate for Importance Sampling	458
13.7.2	Importance Sampling in Bayesian Applications	461
13.8	The EM Algorithm	467
13.8.1	General Background	469
	Exercises	472
Appendix A	Review of Probability	477
A.1	Basic Probability	477
A.2	Mean and Variance	478

A.3	The Normal Distribution	480
A.4	The Mean of a Sample of Random Variables	481
A.5	Sums of Normal Random Variables	482
A.6	The Law of Averages	483
A.7	Higher Moments and the Moment-generating Function	484

Appendix B Probability Distributions 487

B.1	The Bernoulli and Binomial Distributions	487
B.2	The Multinomial Distribution	488
B.3	The Geometric Distribution	490
B.4	The Negative Binomial Distribution	491
B.5	The Hypergeometric Distribution	492
B.6	The Poisson Distribution	493
B.7	The Uniform Distribution	495
B.8	The Exponential Distribution	495
B.9	The Gamma Distribution	497
B.10	The Chi-square Distribution	499
B.11	The Student's t Distribution	502
B.12	The Beta Distribution	504
B.13	The F Distribution	505
	Exercises	507

Appendix C Distributions Quick Reference 509

Solutions to Selected Exercises 513

References 525

Index 531

Preface

Mathematical Statistics with Resampling and R is a one-term undergraduate statistics textbook aimed at sophomores or juniors who have taken a course in probability (at the level of, for instance, Ross (2009), Ghahramani (2004), or Scheaffer and Young (2010)) but may not have had any previous exposure to statistics.

What sets this book apart from other mathematical statistics texts is the use of modern resampling techniques – permutation tests and bootstrapping. We begin with permutation tests and bootstrap methods before introducing classical inference methods. Resampling helps students understand the meaning of sampling distributions, sampling variability, P -values, hypothesis tests, and confidence intervals. We are inspired by the textbooks of Wardrop (1995) and Chance and Rossman (2005), two innovative introductory statistics books that also take a nontraditional approach in the sequencing of topics.

We believe the time is ripe for this book. Many faculty have learned resampling and simulation-based methods in graduate school and/or use them in their own work and are eager to incorporate these ideas into a mathematical statistics course. Students and faculty today have access to computers that are powerful enough to perform resampling quickly.

A major topic of debate about the mathematical statistics course is how much theory to introduce. We want mathematically talented students to get excited about statistics, so we try to strike a balance between theory, computing, and applications. We feel that it is important to demonstrate some rigor in developing some of the statistical ideas presented here, but that mathematical theory should not dominate the text. To keep the size of the text reasonable, we omit some topics such as sufficiency and Fisher information (though we plan to make some omitted topics available as supplements on the text web page <https://sites.google.com/site/ChiharaHesterberg>).

We have compiled the definitions and theorems of the important probability distributions into an appendix (see Appendix B). Instructors who want to prove results on distributional theory can refer to that chapter. Instructors who wish

to skip the theory can continue without interrupting the flow of the statistical discussion.

Incorporating resampling and bootstrapping methods requires that students use statistical software. We use R or RStudio because they are freely available (www.r-project.org or rstudio.com), powerful, flexible, and a valuable tool in future careers. One of us works at Google where there is an explosion in the use of R, with more and more nonstatisticians learning R (the statisticians already know it). We realize that the learning curve for R is high, but believe that the time invested in mastering R is worth the effort. We have written some basic materials on R that are available on the website for this text. We recommend that instructors work through the introductory worksheet with the students on the first or second day of the term in a computer lab if possible.

We had some discussion about whether to include examples of R code using various packages (including **ggplot2**), but we received feedback from colleagues who felt that students should learn to do basic exploratory data analysis and quick diagnostics using base R. And though some R packages exist that implement some of the bootstrap and permutation algorithms that we teach, we felt that students understand and internalize the concepts better if they are required to write the code themselves. We do provide R scripts or R Markdown files with code on our website, and we may include alternate coding using some of the many R packages available.

Statistical computing is necessary in statistical practice and for people working with data in a wide variety of fields. There is an explosion of data more and more data – and new computational methods are continuously being developed to handle this explosion. Statistics is an exciting field; dare we even say sexy?¹

Second Edition: Major changes from the first edition include splitting Chapter 3, the introduction to hypothesis testing, in two. The second half which dealt with categorical data and contingency tables is now its own chapter and moved to later in the textbook. We decided to move mention of the $\alpha = 0.05$ significance level from Chapter 2 to Chapter 8 in an attempt to discourage students' reliance on this threshold. We have also included a new case study using data from Google, plus some discussion on statistical practice. We moved the short chapter on one-way ANOVA available previously on our website into the book; other additions include sections on the bootstrap t test and an introduction to the EM algorithm. We have also added more exercises

1 Try googling “statistics sexy profession.”

and “real” data sets. Throughout the text, we have updated, clarified, or made small changes to the exposition.

Pathways: This textbook contains more than enough material for a one-term undergraduate course. We have written the textbook in such a way that instructors can choose to skip or gloss over some sections if they wish to emphasize others. In some instances, we have labeled a section or subsection with an asterisk (*) to denote it as optional. For classes comprised primarily of students who have no statistics background, a possible sequence includes Chapters 1–10. For courses focused more on applications, instructors could omit, for example, Sections 7.2 and 8.6. For classes in which students come in with an introductory statistics course background, instructors could have students read the first two chapters on their own, beginning the course at Chapter 3. In this case, instructors may wish to spend more time on theory, Bayesian methods, or the topics in Chapter 13, including the parametric bootstrap, the delta method, and importance sampling.

Acknowledgments: This textbook could not have been completed without the assistance of many colleagues and students. In particular, for the first edition, we would like to thank Professor Katherine St. Clair of Carleton College who bravely class tested an early (very!) rough draft in her *Introduction to Statistical Inference* class during winter 2010. In addition, Professor Julie Legler of St. Olaf College adopted the manuscript in her *Statistical Theory* class for fall 2010. Both instructors and students provided valuable feedback that improved the exposition and content of this textbook. For the second edition, we would also like to thank Professors Elaine Newman (Sonoma State University), Eric Nordmoe (Kalamazoo College), and Nick Horton (Amherst College) and Carleton College faculty Katherine St. Clair, Andy Poppick, and Adam Loy for their helpful comments. We thank Ed Lee of Google for the Mobile Ads data and explanation. In addition, we thank Professor Albert Y. Kim (Smith College) who compiled the data sets for the first edition into an R package (**resampled**data).

We would also like to thank Siyuan (Ernest) Liu and Chen (Daisy) Sun, two Carleton College students, for solving many of the exercises in the first edition and writing up the solutions with L^AT_EX.

Finally, the staff at Wiley, including Jon Gurstelle, Amudhapriya Sivamurthy, Kshitija Iyer, Vishnu Narayanan, Kathleen Pagliaro, Steve Quigley, Sanchari Sill, Dean Gonzalez, and Jackie Palmieri provided valuable assistance in preparing both the first and second edition manuscripts for press.

Additional Resources: The authors' web page for this book <https://sites.google.com/site/ChiharaHesterberg> contains R scripts, data sets, tutorials, errata, and supplemental material.

1

Data and Case Studies

Statistics is the art and science of collecting and analyzing data and understanding the nature of variability. Mathematics, especially probability, governs the underlying theory, but statistics is driven by applications to real problems.

In this chapter, we introduce several data sets that we will encounter throughout the text in the examples and exercises.

1.1 Case Study: Flight Delays

If you have ever traveled by air, you probably have experienced the frustration of flight delays. The Bureau of Transportation Statistics maintains data on all aspects of air travel, including flight delays at departure and arrival (<https://www.bts.gov/topics/airlines-and-airports/quick-links-popular-air-carrier-statistics>).

LaGuardia Airport (LGA) is one of three major airports that serves the New York City metropolitan area. In 2008, over 23 million passengers and over 375 000 planes flew in or out of LGA. United Airlines and American Airlines are two major airlines that schedule services at LGA. The data set `FlightDelays` contains information on all 4029 departures of these two airlines from LGA during May and June 2009 (Tables 1.1 and 1.2).

Each row of the data set is an *observation*. Each column represents a *variable* – some characteristic that is obtained for each observation. For instance, on the first observation listed, the flight was a United Airlines plane, flight number 403, destined for Denver, and departing on Friday between 4 and 8 a.m. This data set consists of 4029 observations and 9 variables.

Questions we might ask include the following: Are flight delay times different between the two airlines? Are flight delay times different depending on the day of the week? Are flights scheduled in the morning less likely to be delayed by more than 15 min?

Table 1.1 Partial view of `FlightDelays` data.

Flight	Carrier	FlightNo	Destination	DepartTime	Day
1	UA	403	DEN	4–8 a.m.	Friday
2	UA	405	DEN	8–noon	Friday
3	UA	409	DEN	4–8 p.m.	Friday
4	UA	511	ORD	8–noon	Friday
		⋮			

Table 1.2 Variables in data set `FlightDelays`.

Variable	Description
Carrier	UA=United Airlines, AA=American Airlines
FlightNo	Flight number
Destination	Airport code
DepartTime	Scheduled departure time in 4 h intervals
Day	Day of week
Month	May or June
Delay	Minutes flight delayed (negative indicates early departure)
Delayed30	Departure delayed more than 30 min?
FlightLength	Length of time of flight (minutes)

1.2 Case Study: Birth Weights of Babies

The birth weight of a baby is of interest to health officials since many studies have shown possible links between this weight and conditions in later life, such as obesity or diabetes. Researchers look for possible relationships between the birth weight of a baby and the age of the mother or whether or not she smoked cigarettes or drank alcohol during her pregnancy. The Centers for Disease Control and Prevention (CDC) maintains a database on all babies born in a given year (<http://wonder.cdc.gov/natality-current.html>), incorporating data provided by the US Department of Health and Human Services, the National Center for Health Statistics, and the Division of Vital Statistics. We will investigate different samples taken from the CDC’s database of births.

One data set that we will investigate consists of a random sample of 1009 babies born in North Carolina during 2004 (Table 1.3). The babies in the sample

Table 1.3 Variables in data set `NCBirths2004`.

Variable	Description
Age	Mother's age
Tobacco	Mother used tobacco?
Gender	Gender of baby
Weight	Weight at birth (grams)
Gestation	Gestation time (weeks)

had a gestation period of at least 37 weeks and were single births (i.e. not a twin or triplet).

In addition, we will also investigate a data set, `Girls2004`, consisting of a random sample of 40 baby girls born in Alaska and 40 baby girls born in Wyoming. These babies also had a gestation period of at least 37 weeks and were single births.

The data set `TXBirths2004` contains a random sample of 1587 babies born in Texas in 2004. In this case, the sample was not restricted to single births, nor to a gestation period of at least 37 weeks. The numeric variable `Number` indicates whether the baby was a single birth, or one of a twin, triplet, and so on. The variable `Multiple` is a factor variable indicating whether or not the baby was a multiple birth.

1.3 Case Study: Verizon Repair Times

Verizon is the primary local telephone company (incumbent local exchange carrier (ILEC)) for a large area of the Eastern United States. As such, it is responsible for providing repair service for the customers of other telephone companies known as competing local exchange carriers (CLECs) in this region. Verizon is subject to fines if the repair times (the time it takes to fix a problem) for CLEC customers are substantially worse than those for Verizon customers.

The data set `Verizon` contains a sample of repair times for 1664 ILEC and 23 CLEC customers (Table 1.4). The mean repair times are 8.4 h for ILEC

Table 1.4 Variables in data set `Verizon`.

Variable	Description
Time	Repair times (in hours)
Group	ILEC or CLEC

customers and 16.5 h for CLEC customers. Could a difference this large be easily explained by chance?

1.4 Case Study: Iowa Recidivism

When a person is released from prison, will he or she relapse into criminal behavior and be sent back? The state of Iowa tracks offenders over a 3-year period and records the number of days until recidivism for those who are readmitted to prison. The Department of Corrections uses this recidivism data to determine whether or not their strategies for preventing offenders from relapsing into criminal behavior are effective.

The data set `Recidivism` contains all offenders convicted of either a misdemeanor or felony who were released from an Iowa prison during the 2010 fiscal year (ending in June) (Table 1.5). There were 17 022 people released in that period, of whom 5386 were sent back to prison in the following 3 years (through the end of the 2013 fiscal year).¹

The recidivism rate for those under the age of 25 years was 36.5% compared with 30.6% for those 25 years or older. Does this indicate a real difference in the behavior of those in these age groups, or could this be explained by chance variability?

Table 1.5 Variables in data set `Iowa Recidivism`.

Variable	Description
Gender	F, M
Race	American Indian or Alaska Native Hispanic, American Indian or Alaska Native Non-Hispanic, Asian or Pacific Islander Hispanic, Asian or Pacific Islander NonHispanic, Black, Black Hispanic, Black Non-Hispanic, White, White Hispanic, White Non-Hispanic
Age	Age at release: under 25, 25–34, 35–44, 45–54, and 55 and older
Age25	Under 25, over 25 (binary)
Offense	Original conviction: felony or misdemeanor
Recid	Recidivate? No, yes
Type	New (crime), no recidivism, tech (technical violation, such as a parole violation)
Days	Number of days to recidivism; NA if no recidivism

¹ (<https://data.iowa.gov/Public-Safety/3-Year-Recidivism-for-Offenders-Released-from-Pris/mw8r-vqy4>).

1.5 Sampling

In analyzing data, we need to determine whether the data represent a *population* or a *sample*. A *population* represents all the individual cases, whether they are babies, fish, cars, or coin flips. The data from flight delays case study in Section 1.1 are *all* the flight departures of United Airlines and American Airlines out of LGA in May and June 2009; thus, this data set represents the population of all such flights. On the other hand, the North Carolina data set contains only a subset of 1009 births from over 100 000 births in North Carolina in 2004. In this case, we will want to know how representative statistics computed from this sample are for the entire population of North Carolina babies born in 2004.

Populations may be finite, such as births in 2004, or infinite, such as coin flips or births next year.

Throughout this book, we will talk about drawing random samples from a population. We will use capital letters (e.g. X , Y , Z , and so on) to denote random variables and lower-case letters (e.g. x_1 , x_2 , x_3 , and so on) to denote actual values or data.

There are many kinds of random samples. Strictly speaking, a “random sample” is any sample obtained using a random procedure. However, in this book we use *random sample* to mean a sample of independent and identically distributed (i.i.d.) observations from the population, if the population is infinite.

For instance, suppose you toss a fair coin 20 times and consider each head a “success.” Then your sample consists of the random variables X_1, X_2, \dots, X_{20} , each a Bernoulli random variable with success probability $1/2$. We use the notation $X_i \sim \text{Bern}(1/2)$, $i = 1, 2, \dots, 20$.

If the population of interest is finite $\{x_1, x_2, \dots, x_N\}$, we can choose a random sample as follows: Label N balls with the numbers $1, 2, \dots, N$ and place them in an urn. Draw a ball at random, record its value $X_1 = x_{i_1}$, and then replace the ball. Draw another ball at random, record its value, $X_2 = x_{i_2}$, and replace. Continue until you have a sample $x_{i_1}, x_{i_2}, \dots, x_{i_n}$. This is *sampling with replacement*. For instance, if $N = 5$ and $n = 2$, then there are $5 \times 5 = 25$ different samples of size 2 (where order matters). (Note: By “order matters” we do not imply that order matters in practice, rather we mean that we keep track of the order of the elements when enumerating samples. For instance, the set $\{a, b\}$ is different from $\{b, a\}$.)

However, in most real situations, for example, in conducting surveys, we do not want to have the same person polled twice. So we would sample *without replacement*, in which case, we will not have independence. For instance, if you wish to draw a sample of size $n = 2$ from a population of $N = 10$ people, then the probability of any one person being selected is $1/10$. However, after having chosen that first person, the probability of any one of the remaining people being chosen is now $1/9$.

In cases where populations are very large compared to the sample size, calculations under sampling without replacement are reasonably approximated by calculations under sampling with replacement.

Example 1.1 Consider a population of 1000 people, 350 of whom are smokers, and the rest are nonsmokers. If you select 10 people at random but with replacement, then the probability that 4 are smokers is $\binom{10}{4} (350/1000)^4 (650/1000)^6 \approx 0.2377$. If you select without replacement, then the probability is $\binom{350}{4} \binom{650}{6} / \binom{1000}{10} \approx 0.2388$. \square

1.6 Parameters and Statistics

When discussing numeric information, we will want to distinguish between populations and samples.

Definition 1.1 A *parameter* is a (numerical) characteristic of a population or of a probability distribution.

A *statistic* is a (numerical) characteristic of data. \parallel

Any function of a parameter is also a parameter; any function of a statistic is also a statistic. When the statistic is computed from a random sample, it is itself random, and hence is a random variable.

Example 1.2 μ and σ are parameters of the normal distribution with pdf $f(x) = (1/\sqrt{2\pi}\sigma)e^{-(x-\mu)^2/(2\sigma^2)}$.

The variance σ^2 and *signal-to-noise ratio* μ/σ are also parameters. \square

Example 1.3 If X_1, X_2, \dots, X_n are a random sample, then the mean $\bar{X} = 1/n \sum_{i=1}^n X_i$ is a statistic. \square

Example 1.4 Consider the population of all babies born in the United States in 2017. Let μ denote the average weight of all these babies. Then μ is a parameter. The average weight of a sample of 2500 babies born in that year is a statistic. \square

Example 1.5 If we consider the population of all adults in the United States today, the proportion p who approve of the president's job performance is a parameter. The fraction \hat{p} who approve in any given sample is a statistic. \square

Example 1.6 The average weight of 1009 babies in the North Carolina case study in Section 1.2 is 3448.26 g. This average is a statistic. \square

Example 1.7 If we survey 1000 adults and find that 60% intend to vote in the next presidential election, then $\hat{p} = 0.60$ is a statistic: It estimates the parameter p , the proportion of all adults who intend to vote in the next election. \square

1.7 Case Study: General Social Survey

The General Social Survey (GSS) is a major survey that has tracked American demographics, characteristics, and views on social and cultural issues since the 1970s. It is conducted by the National Opinion Research Center (NORC) at the University of Chicago. Trained interviewers meet face to face with the adults chosen for the survey and question them for about 90 min in their homes.

The GSS case study includes the responses of 2765 participants selected in 2002 to about a dozen questions, listed in Table 1.6. For example, one of the questions (`SpendEduc`) asked whether the respondent believed that the

Table 1.6 Variables in data set `GSS2002`.

Variable	Description
Region	Interview location
Gender	Gender of respondent
Race	Race of respondent: White, Black, Other
Marital	Marital status
Education	Highest level of education
Happy	General happiness
Income	Respondent's income
PolParty	Political party
Politics	Political views
Marijuana	Legalize marijuana?
DeathPenalty	Death penalty for murder?
OwnGun	Have gun at home?
GunLaw	Require permit to buy a gun?
SpendMilitary	Amount government spends on military
SpendEduc	Amount government spends on education
SpendEnv	Amount government spends on the environment
SpendSci	Amount government spends on science
Pres00	Whom did you vote for in the 2000 presidential election?
Postlife	Believe in life after death?

amount of money being spent on the nation's education system was too little, too much, or the right amount.

We will analyze the GSS data to investigate questions such as the following: Is there a relationship between the gender of an individual and whom they voted for in the 2000 presidential election? Are people who live in certain regions happier? Are there educational differences in support for the death penalty? These data are archived at the Computer-assisted Survey Methods Program at the University of California (www.sda.berkeley.edu).

1.8 Sample Surveys

“Who do you plan to vote for in the next presidential election?” “Would you purchase our product again in the future?” “Do you smoke cigarettes? If yes, how old were you when you first started?” Questions such as these are typical of sample surveys. Researchers want to know something about a population of individuals, whether they are registered voters, online shoppers, or American teenagers, but to poll every individual in the population – that is, to take a *census* – is impractical and costly. Thus, researchers will settle for a sample from the target population. But if, say, 60% of those in your sample of 1000 adults intend to vote for candidate Wong in the next election, how close is this to the actual percentage who will vote for Wong? How can we be sure that this sample is truly representative of the population of all voters? We will learn techniques for *statistical inference*, drawing a conclusion about a population based on information about a sample.

When conducting a survey, researchers will start with a *sampling frame* – a list from which the researchers will choose their sample. For example, to survey all students at a college, the campus directory listing could be a sampling frame. For pre-election surveys, many polling organizations use a sampling frame of registered voters. Note that the choice of sampling frame could introduce the problem of *undercoverage*: omitting people from the target population in the survey. For instance, young people were missed in many pre-election surveys during the 2008 Obama–McCain presidential race because they had not yet registered to vote.

Once the researchers have a sampling frame, they will then draw a random sample from this frame. Researchers will use some type of *probability (scientific) sampling scheme*, that is, a scheme that gives everybody in the population a positive chance of being selected. For example, to obtain a sample of size 10 from a population of 100 individuals, write each person's name on a slip of paper, put the slips of paper into a basket, and then draw out 10 slips of paper. Nowadays, statistical software is used to draw random samples from a sampling frame.

Another basic survey design uses *stratified sampling*: The population is divided into nonoverlapping strata, and then random samples are drawn from each stratum. The idea is to group individuals who are similar in some characteristic into homogeneous groups, thus reducing variability. For instance, in a survey of university students, a researcher might divide the students by class: first year, sophomores, juniors, seniors, and graduate students. A market analyst for an electronics store might choose to stratify customers based on income levels.

In *cluster sampling*, the population is divided into nonoverlapping clusters, and then a random sample of clusters is drawn. Every person in a chosen cluster is then interviewed for the survey. An airport wanting to conduct a customer satisfaction survey might use a sampling frame of all flights scheduled to depart from the airport on a certain day. A random sample of flights (clusters) is chosen, and then all passengers on these flights are surveyed. A modification of this design might involve sampling in stages: For instance, the analysts might first choose a random sample of flights, and then from each flight choose a random sample of passengers.

The GSS uses a more complex sampling scheme in which the sampling frame is a list of counties and county equivalents (standard metropolitan statistical areas) in the United States. These counties are stratified by region, age, and race. Once a sample of counties is obtained, a sample of block groups and enumeration districts is selected, stratifying these by race and income. The next stage is to randomly select blocks and then interview a specific number of men and women who live within these blocks.

Indeed, all major polling organizations such as Gallup or Roper as well as the GSS use a *multistage* sampling design. In this book, we use the GSS data or polling results for examples as if the survey design used simple random sampling. Calculations for more complex sampling scheme are beyond the scope of this book, and we refer the interested reader to Lohr (1991) for details.

1.9 Case Study: Beer and Hot Wings

Carleton student Nicki Catchpole conducted a study of hot wings and beer consumption at the Williams Bar in the Uptown area of Minneapolis (N. Catchpole, private communication). She asked patrons at the bar to record their consumption of hot wings and beer over the course of several hours. She wanted to know if people who ate more hot wings would then drink more beer. In addition, she investigated whether or not gender had an impact on hot wings or beer consumption.

The data for this study are in *Beerwings* (Table 1.7). There are 30 observations and 3 variables.

Table 1.7 Variables in data set *Beerwings*.

Variable	Description
Gender	Male or female
Beer	Ounces of beer consumed
Hot Wings	Number of hot wings eaten

1.10 Case Study: Black Spruce Seedlings

Black spruce (*Picea mariana*) is a species of a slow-growing coniferous tree found across the northern part of North America. It is commonly found on wet organic soils. In a study conducted in the 1990s, a biologist interested in factors affecting the growth of the black spruce planted its seedlings on sites located in boreal peatlands in northern Manitoba, Canada (Camil et al. (2010)).

The data set *Spruce* contains a part of the data from the study (Table 1.8). Seventy-two black spruce seedlings were planted in four plots under varying conditions (fertilizer–no fertilizer, competition–no competition), and their heights and diameters were measured over the course of 5 years.

The researcher wanted to see whether the addition of fertilizer or the removal of competition from other plants (by weeding) affected the growth of these seedlings.

1.11 Studies

Researchers carry out studies to understand the conditions and causes of certain outcomes: Does smoking cause lung cancer? Do teenagers who smoke marijuana tend to move on to harder drugs? Do males eat more hot wings than females? Do black spruce seedlings grow taller in fertilized plots?

Table 1.8 Variables in data set *Spruce*.

Variable	Description
Tree	Tree number
Competition	C (competition), CR (competition removed)
Fertilizer	F (fertilized), NF (not fertilized)
Height0	Height (cm) of seedling at planting
Height5	Height (cm) of seedling at year 5
Diameter0	Diameter (cm) of seedling at planting
Diameter5	Diameter (cm) of seedling at year 5
Ht.change	Change (cm) in height
Di.change	Change (cm) in diameter

The beer and hot wings case study in Section 1.9 is an example of an *observational study*, a study in which researchers observe participants but do not influence the outcome. In this case, the student just recorded the number of hot wings eaten and beer consumed by the patrons of Williams Bar.

Example 1.8 The first Nurses' Health Study is a major observational study funded by the National Institutes of Health. Over 12 000 registered female nurses who, in 1976, were married, between the ages of 33 and 55 years, and who lived in the 11 most populous states have been responding every 2 years to written questions about their health and lifestyle, including smoking habits, hormone use, and menopause status. Many results on women's health have come out of this study, such as finding an association between taking estrogen after menopause and lowering the risk of heart disease, and determining that for nonsmokers there is no link between taking birth control pills and developing heart disease.

Because this is an observational study, no *cause-and-effect* conclusions can be drawn. For instance, we cannot state that taking estrogen after menopause will *cause* a lowering of the risk for heart disease. In an observational study, there may be many unrecorded or hidden factors that impact the outcomes. Also, because the participants in this study are registered nurses, we need to be careful about making inferences about the general female population. Nurses are more educated and more aware of health issues than the average person. □

On the other hand, the black spruce case study in Section 1.10 was an *experiment*. In an experiment, researchers will manipulate the environment in some way to observe the response of the objects of interest (people, mice, ball bearings, etc.). When the objects of interest in an experiment are people, we refer to them as *subjects*; otherwise, we call them *experimental units*. In this case, the biologist randomly assigned the experimental units – the seedlings – to plots subject to four *treatments*: fertilization with competition, fertilization without competition, no fertilization with competition, and no fertilization with no competition. He then recorded their height over a period of several years.

A key feature in this experiment was the *random assignment* of the seedlings to the treatments. The idea is to spread out the effects of unknown or uncontrollable factors that might introduce unwanted variability into the results. For instance, if the biologist had planted all the seedlings obtained from one particular nursery in the fertilized, no competition plot and subsequently recorded that these seedlings grew the least, then he would not be able to discern whether this was due to this particular treatment or due to some possible problem with seedlings from this nursery. With random assignment of treatments, the seedlings from this particular nursery would usually be spread out over the four treatments. Thus, the differences between the treatment groups should be due to the treatments (or chance).

Example 1.9 Knee osteoarthritis (OA) that results in deterioration of cartilage in the joint is a common source of pain and disability for the elderly population. In a 2008 paper, “Tai Chi is effective in treating knee osteoarthritis: A randomized controlled trial,” Wang et al. (2009) at Tufts University Medical School describe an experiment they conducted to see whether practicing tai chi, a style of Chinese martial arts, could alleviate pain from OA. Forty patients over the age of 65 with confirmed knee OA but otherwise in good health were recruited from the Boston area. Twenty were randomly assigned to attend twice weekly 60 min sessions of tai chi for 12 weeks. The remaining 20 participants, the *control group*, attended twice weekly 60 min sessions of instructions on health and nutrition, as well as some stretching exercises.

At the end of the 12 weeks, those in the tai chi group reported a significant decrease in knee pain. Because the subjects were randomly assigned to the two treatments, the researchers can assert that the tai chi sessions lead to decrease in knee pain due to OA. Note that because the subjects were recruited, we need to be careful about making an inference about the general elderly population: People who voluntarily sign up to be in an experiment may be different from other people. □

Another important feature of a well-designed experiment is *blinding*: A *double-blind* experiment is one in which neither the researcher nor the subject knows who is receiving which treatment. An experiment is *single-blinded* if just the researcher or the subject (but not both) knows who is receiving which treatment. Blinding is important in reducing *bias*, the systematic favoring of one outcome over another.

For instance, suppose in a clinical trial to test the efficacy of a new drug for a disease, the subjects know whether they are receiving the drug or a placebo (a pill or drug with no therapeutic effect). Those on the placebo might feel that the trial is a waste of time and drop out, or perhaps seek additional treatment elsewhere. On the other hand, if the researcher knows that a subject received the drug, he or she might behave differently toward the subject, perhaps by asking leading questions that result in responses that appear to suggest relief from the disease.

1.12 Google Interview Question: Mobile Ads Optimization

The following question was posted on an internal Google statistics email list:

I have a pre v post comparison I'm trying to make where alternative hypothesis is $\text{pre.mean.error} > \text{post.mean.error}$. My distribution for these samples are both right skewed as shown below. Anyone know what test method would be best suited for this type of situation?