

Introduction to Probability and Statistics for Data Science

Introduction to Probability and Statistics for Data Science provides a solid course in the fundamental concepts, methods, and theory of statistics for students in statistics, data science, biostatistics, engineering, and physical science programs. It teaches students to understand, use, and build on modern statistical techniques for complex problems. The authors develop the methods from both an intuitive and mathematical angle, illustrating with simple examples how and why the methods work. More complicated examples, many of which incorporate data and code in R, show how the method is used in practice. Through this guidance, students get the big picture about how statistics works and can be applied. This text covers more modern topics such as regression trees, large-scale hypothesis testing, bootstrapping, MCMC, time series, and fewer theoretical topics such as the Cramer–Rao lower bound and the Rao–Blackwell theorem. It features more than 250 high-quality figures, 180 of which involve actual data. Data and R code are available on the book’s website so that students can reproduce the examples and complete hands-on exercises.

Steven E. Rigdon is Professor of Biostatistics at Saint Louis University. He is a fellow of the American Statistical Association and is the author of *Statistical Methods for the Reliability of Repairable Systems*, *Calculus*, 8th and 9th editions, *Monitoring the Health of Populations by Tracking Disease Outbreaks* (2020), and *Design of Experiments for Reliability Achievement* (2022). He has received the Waldo Vizeau Award for technical contributions to quality, the Soren Bisgaard Award, and the Paul Simon Award for linking teaching and research. He is also Distinguished Research Professor Emeritus at Southern Illinois University Edwardsville.

Ronald D. Fricker, Jr. is Vice Provost for Faculty Affairs at Virginia Tech, where he has served as head of the Department of Statistics, Senior Associate Dean in the College of Science, and, subsequently, interim dean of the college. He is the author of *Introduction to Statistical Methods for Biosurveillance* (2013) and, with Steve Rigdon, *Monitoring the Health of Populations by Tracking Disease Outbreaks* (2020). He is a fellow of the American Statistical Association, a fellow of the American Association for the Advancement of Science, and an elected member of the Virginia Academy of Science, Engineering, and Medicine.

Douglas C. Montgomery is Regents’ Professor and ASU Foundation Professor of Engineering at Arizona State University. He is an Honorary Member of the American Society for Quality, a fellow of the American Statistical Association, a fellow of the Institute of Industrial and Systems Engineering, and a fellow of the Royal Statistical Society. He is the author of 15 other books, including *Design and Analysis of Experiments*, 10th edition (2020) and *Design of Experiments for Reliability Achievement* (2022). He has received the Shewhart Medal, the Distinguished Service Medal, and the Brumbaugh Award from the ASQ, the Deming Lecture Award from the ASA, the Greenfield Medal from the Royal Statistical Society, and the George Box Medal from the European Network for Business and Industrial Statistics.

“This book serves as an excellent resource for students with diverse backgrounds, offering a thorough exploration of fundamental topics in statistics. The clear explanation of concepts, methods, and theory, coupled with an abundance of practical examples, provides a solid foundation to help students understand statistical principles and bridge the gap between theory and application. This book offers invaluable insights and guidance for anyone seeking to master the principles of statistics. I highly recommend adopting this book for my future statistics class.”

Haijun Gong, *Saint Louis University*

“Professors Rigdon, Fricker and Montgomery have put together an impressive volume that covers not only basic probability and basic statistics, but also includes extensions in a number of directions, all of which have immediate relevance to the work of practitioners in quantitative fields. Suffused with common sense and insights about real data and problems, it is both approachable and precise. I’m excited about the inclusion of material on power and on multiple testing, both of which will help users become smarter about what their analyses can do, and I applaud their omission of too much theory. I also appreciate their use of R and of real data. This would be an excellent text for undergraduate or graduate-level data analysts.”

Sam Buttrey, *Naval Postgraduate School (NPS)*

“This is a comprehensive and rich book that extends foundational concepts in statistics and probability in easily accessible form into data science as an integrated discipline. The reader applies and validates theoretical concepts in R and connects results from R back to the theory across many methods: from descriptive statistics to Bayesian models, time series, generalized linear models and more. Thoroughly enjoyable!”

Oliver Schabenberger, *Virginia Tech Academy of Data Science*

Introduction to Probability and Statistics for Data Science

with R

Steven E. Rigdon

Saint Louis University

Ronald D. Fricker, Jr.

Virginia Polytechnic Institute and State University

Douglas C. Montgomery

Arizona State University



Shaftesbury Road, Cambridge CB2 8EA, United Kingdom
One Liberty Plaza, 20th Floor, New York, NY 10006, USA
477 Williamstown Road, Port Melbourne, VIC 3207, Australia
314-321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi - 110025, India
103 Penang Road, #05-06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of Cambridge University Press & Assessment,
a department of the University of Cambridge

We share the University's mission to contribute to society through the pursuit of
education, learning and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/highereducation/isbn/9781107113046

DOI: 10.1017/9781316286166

© Steven E. Rigdon, Ronald D. Fricker, Jr., and Douglas C. Montgomery 2025

This publication is in copyright. Subject to statutory exception and to the provisions
of relevant collective licensing agreements, no reproduction of any part may take place
without the written permission of Cambridge University Press & Assessment.

When citing this work, please include a reference to the DOI 10.1017/9781316286166

First published 2025

Printed in Mexico by Litográfica Ingramex, S.A. de C.V.

A catalogue record for this publication is available from the British Library.

A Cataloging-in-Publication data record for this book is available from the Library of Congress.

ISBN 978-1-107-11304-6 Hardback

ISBN 978-1-009-56835-7 Paperback

Additional resources for this publication at www.cambridge.org/ProbStatsforDS.

Cambridge University Press has no responsibility for the persistence or accuracy of
URLs for external or third-party internet websites referred to in this publication
and does not guarantee that any content on such websites is, or will remain,
accurate or appropriate.

Steve Rigdon

To my wife Pat, who has always supported me and been at my side.

Ron Fricker

To my spouse, Christine: *Tu ventus sub alis meis es.*

And to my first statistics professor, Randy Spoeri: You introduced me to the subject and made it fun.

Doug Montgomery

To Cheryl, who has always supported and encouraged me. And to the memory of my first statistics professor, Ray Myers, mentor, colleague, collaborator, and friend.



Contents

Preface	xiii		
1 Introduction	1		
1.1 Data Science and Statistics	2		
1.2 More on Statistics	3		
1.2.1 Populations and Samples	4		
1.2.2 Descriptive versus Inferential Statistics	5		
1.3 An Introduction to R	6		
1.4 Descriptive Statistics	7		
1.4.1 Types of Data	8		
1.4.2 Example Data: US Domestic Flights from 1987 to 2008	9		
1.5 Cross-Sectional Data	9		
1.5.1 Measures of Location	9		
1.5.2 Measures of Variation	16		
1.5.3 Measures of How Two Variables Co-vary	18		
1.5.4 Other Summary Statistics	23		
1.6 Tabular Summaries of Data	25		
1.7 Chapter Summary	27		
1.8 Problems	29		
2 Data Visualization	31		
2.1 Introduction	31		
2.2 Traditional Statistical Graphics	32		
2.2.1 Bar Charts	32		
2.2.2 Pie Charts	35		
2.2.3 Histograms	35		
2.2.4 Lattice (or Trellis) Plots	38		
2.2.5 Box Plots	40		
2.2.6 Scatterplots	43		
2.3 Graphics for Longitudinal Data	46		
2.3.1 Time Series Plots	47		
2.3.2 Repeated Cross-Sectional Plots	48		
2.3.3 Autocorrelation Plots	50		
2.4 Chapter Summary	50		
2.5 Problems	52		
3 Basic Probability	54		
3.1 Introduction	54		
3.2 Events and Sample Spaces	54		
3.2.1 Probability Axioms	56		
3.2.2 Union of Events	57		
3.2.3 Intersection of Independent Events	59		
3.2.4 Complementary Events	60		
3.2.5 Conditional Probability	62		
3.3 Calculating Probabilities	63		
3.3.1 Sample Point Method	64		
3.3.2 Counting Sample Points	66		
3.3.3 Combining Events	70		
3.4 Bringing It All Together	71		
3.4.1 Law of Total Probability	71		
3.4.2 Bayes' Theorem	74		
3.5 Chapter Summary	77		
3.6 Problems	78		

4	Random Variables	82	6.5	Gamma and Weibull Distributions	189
4.1	Introduction	82	6.5.1	Gamma Distribution	189
4.2	Discrete Random Variables	82	6.5.2	Weibull Distribution	192
4.2.1	Probability Mass Function	82	6.6	Distributions Related to the Normal	196
4.2.2	Cumulative Distribution Function	86	6.6.1	Chi-square (χ^2) Distribution	196
4.2.3	Expected Value	88	6.6.2	t Distribution	198
4.2.4	Variance and Standard Deviation	91	6.6.3	F Distribution	201
4.3	Continuous Random Variables	93	6.7	Beta Distribution	202
4.3.1	Probability Density Function	93	6.8	Transformations	205
4.3.2	Cumulative Distribution Function	97	6.8.1	Simulating from Distributions	207
4.3.3	Expected Value	101	6.9	Moment Generating Functions	209
4.3.4	Variance and Standard Deviation	101	6.10	Quantile–Quantile Plots	215
4.4	Expected Value and Variance Properties	102	6.11	Chapter Summary	220
4.5	Joint Distributions for Discrete Random Variables	105	6.12	Problems	221
4.6	Conditional Distributions for Discrete Random Variables	111	7	About Data and Data Collection	226
4.7	Joint Distributions for Continuous Random Variables	113	7.1	Introduction	226
4.8	Conditional Distributions for Continuous Random Variables	121	7.2	Data and the Scientific Method	228
4.9	Conditioning on a Random Variable	123	7.3	Experimental vs. Observational Data	231
4.10	Chapter Summary	126	7.3.1	Convenience vs. Probability Sampling	234
4.11	Problems	127	7.4	Accuracy vs. Precision	234
5	Discrete Distributions	132	7.5	Types of Random Samples	236
5.1	Introduction	132	7.5.1	Sources of Bias	237
5.2	Binomial Distribution	132	7.6	Types of Error	239
5.3	Geometric Distribution	141	7.7	Historical Gaffes in Data Collection	240
5.4	Negative Binomial	146	7.8	Chapter Summary	241
5.5	Hypergeometric Distribution	149	7.9	Problems	242
5.6	Poisson Distribution	154	8	Sampling Distributions	244
5.7	Multinomial Distribution	160	8.1	Introduction	244
5.8	Chapter Summary	164	8.2	Linear Combinations of Random Variables	244
5.9	Problems	166	8.3	Sampling Distributions for Sums and Means	249
6	Continuous Distributions	170	8.4	Sampling Distribution for the Sample Variance	253
6.1	Introduction	170	8.5	The Central Limit Theorem	255
6.2	Uniform Distribution	170	8.6	Normal Approximation to the Binomial	259
6.3	Exponential Distribution	173	8.7	Sampling Distributions for Proportions	263
6.4	Normal Distribution	180	8.8	Tchebysheff's Theorem and the Law of Large Numbers	265
6.4.1	Standardizing	183			
6.4.2	Bivariate and Multivariate Normal Distributions	186			

8.9	Chapter Summary	268	11.5	Testing the Mean: Variance Known	357
8.10	Problems	269	11.5.1	Hypothesis Tests for the Mean from a Population with Known Variance	357
9	Point Estimation	273	11.5.2	Power	360
9.1	Introduction and Intuitive Estimators	273	11.6	Testing the Mean: Variance Unknown	364
9.2	Estimation Criteria	275	11.7	Testing a Proportion	370
9.2.1	Unbiased Estimators	275	11.8	Testing the Variance	373
9.2.2	Consistent Estimators	277	11.9	Likelihood Ratio Tests	374
9.3	Method of Moments	279	11.10	Chapter Summary	380
9.4	Maximum Likelihood	283	11.11	Problems	381
9.5	Approximating MLEs	289	12	Hypothesis Tests for Two or More Populations	386
9.6	Sufficiency	294	12.1	Introduction	386
9.7	Chapter Summary	298	12.2	Testing Two Independent Samples	386
9.8	Problems	299	12.2.1	Comparing Two Means	386
10	Confidence Intervals	302	12.2.2	Comparing Two Proportions	398
10.1	Introduction	302	12.2.3	Comparing Variances	402
10.2	Basic Properties	303	12.3	Testing Paired Samples	404
10.3	Large Sample Confidence Intervals	307	12.4	Single-Factor Analysis of Variance	409
10.4	Small Sample Confidence Intervals	312	12.5	Two-Factor ANOVA	425
10.5	Confidence Intervals for Differences	316	12.6	Other Designs for Experiments	430
10.5.1	Confidence Intervals for Differences of Proportions	317	12.6.1	Two-Level Factorial Designs	431
10.5.2	Confidence Intervals for Differences in Means	319	12.6.2	Fractional Factorial Designs	436
10.5.3	Confidence Interval for Paired Data	324	12.6.3	Block Designs	439
10.6	Determining the Sample Size	326	12.6.4	Some Experimental Design Principles	442
10.7	Confidence Intervals from Complex Survey Data	330	12.7	Power	444
10.7.1	Sampling from a Finite Population	330	12.7.1	Power for Two-Sample t -Test	445
10.7.2	Stratified Random Samples	333	12.7.2	Power for One-Factor ANOVA	448
10.7.3	Cluster Sampling	338	12.8	Chapter Summary	451
10.7.4	Secondary Data Sources	339	12.9	Problems	453
10.7.5	Software for Analyzing Data from Complex Surveys	341	13	Hypothesis Tests for Categorical Data	459
10.8	Chapter Summary	343	13.1	Introduction	459
10.9	Problems	344	13.2	Goodness-of-Fit Tests	460
11	Hypothesis Testing	348	13.3	Contingency Tables: Testing Independence	465
11.1	Introduction	348	13.4	Contingency Tables: Homogeneity	471
11.2	Elements of a Statistical Test	350	13.5	Fisher's Exact Test	473
11.3	Power	354			
11.4	P -values	356			

13.6	The Continuity Correction and Simulation	479	15.6	Noninformative Priors	596
13.7	McNemar's Test	481	15.7	Simulation Methods	599
13.8	Higher-Dimensional Tables and Simpson's Paradox	485	15.7.1	Metropolis–Hastings Algorithm	601
13.9	Chapter Summary	488	15.7.2	The Gibbs Sampling Algorithm	604
13.10	Problems	489	15.8	Hierarchical Bayes Models	606
14	Regression	493	15.9	Chapter Summary	613
14.1	Introduction	493	15.10	Problems	614
14.1.1	Prediction vs. Explanation	493	16	Time Series Methods	618
14.1.2	Terminology	495	16.1	Introduction	618
14.1.3	A Working Example	495	16.2	Using R for Time Series	622
14.2	Simple Linear Regression	496	16.3	Numerical Description of Time Series	623
14.3	Properties of the Least Squares Estimators	503	16.4	Exponential Smoothing Methods	628
14.4	Inference for Parameters of the Simple Linear Regression Model	510	16.5	Autoregressive Integrated Moving Average (ARIMA) Models	635
14.5	Matrix Formulation of Simple Linear Regression	516	16.6	Chapter Summary	640
14.6	Joint Confidence Regions	518	16.7	Problems	642
14.7	Confidence and Prediction Intervals for Responses	520	17	Estimating the Standard Error: Analytic Approximations, the Jackknife, and the Bootstrap	645
14.8	Optimal Selection of Levels of Predictor Variables	526	17.1	Introduction	645
14.9	The ANOVA Table for Simple Linear Regression	528	17.2	Analytic Approximations to the Standard Error of an Estimator	646
14.10	Linear Models in More than One Predictor	531	17.3	The Jackknife	656
14.11	Indicator Variables	539	17.4	The Bootstrap	662
14.12	Polynomial and Nonlinear Regression	542	17.4.1	Bootstrap Confidence Intervals Based on Percentiles	667
14.13	Inference for a Linear Combination of Model Parameters	546	17.5	Parametric Bootstrap	670
14.14	Correlation	551	17.6	Bootstrapping in R	674
14.15	R^2 and Adjusted R^2	561	17.7	Chapter Summary	681
14.16	Model Checking	565	17.8	Problems	682
14.16.1	Normal Probability Plots	566	18	Generalized Linear Models and Regression Trees	684
14.16.2	Plot of Residuals against the Fitted Values	567	18.1	Logistic Regression	684
14.17	Chapter Summary	568	18.2	Multinomial Logistic Regression	698
14.18	Problems	570	18.3	Poisson Regression	703
15	Bayesian Methods	574	18.4	Generalized Linear Models	706
15.1	Introduction	574	18.5	Regression Trees	707
15.2	Bayes' Theorem	575	18.6	Discrimination and Classification	714
15.3	The Bayesian Paradigm	581			
15.4	Two Paradoxes	585			
15.5	Conjugate Priors	588			

18.6.1	$K = 2$ Groups and $p = 1$ Variable	718	19.5	Cross-Validation with Classification Data	765
18.6.2	K Groups and $p = 1$ Variable	724	19.6	Chapter Summary	767
18.6.3	K Groups and p Variables	725	19.7	Problems	769
18.6.4	Quadratic Discriminant Analysis	737	20	Large-Scale Hypothesis Testing	773
18.6.5	Dealing with Estimated Parameters	739	20.1	Review of Hypothesis Testing	773
18.6.6	Choosing Between Linear and Quadratic Discriminant Analysis	739	20.2	Testing Multiple Hypotheses	775
18.7	Logistic Regression for Classification	740	20.3	The FWER and the Bonferroni Correction	780
18.8	Chapter Summary	745	20.4	Holm's Method	783
18.9	Problems	746	20.5	The False Discovery Rate	785
19	Cross-Validation and Estimates of Prediction Error	751	20.6	Simultaneous Confidence Intervals	789
19.1	Overfitting and Underfitting	751	20.7	Tukey's Method	791
19.2	Cross-Validation	754	20.8	Scheffé's Method	794
19.2.1	Splitting the Data at Random	755	20.9	Chapter Summary	801
19.3	Leave-One-Out Cross-Validation	760	20.10	Problems	802
19.4	k -Fold Cross-Validation	762	References	805	
			Index	809	



This book is designed for students in statistics, data science, biostatistics, engineering, and mathematics programs who need a solid course in the fundamental concepts, methods, and theory of statistics. Our goal is to give students enough background in the methods and theory of statistics that they can understand modern techniques used in statistics and be able to apply them in the practice of data science.

We had to make some difficult choices regarding topic coverage. We do cover the important concepts of statistics, including maximum likelihood, the information matrix, power, etc., because these are needed for a student to be a successful statistician. When we cover maximum likelihood estimation, we specifically cover the method of approximating the maximum of the (log) likelihood function. Nowadays, data are so plentiful that we are often faced with testing multiple null hypotheses. Holm's method and the Benjamini–Hochberg method are derived and applied to real problems. There are a number of statistical methods that were developed in the late twentieth and early twenty-first centuries, including regression trees, large-scale hypothesis testing, methods of cross-validation, the bootstrap, Markov chain Monte Carlo, and others. We address the optimal selection of levels of a predictor variable to maximize the information we obtain; this leads to an introduction to the topic of optimal design. With some exceptions, these techniques have not found their way into introductory textbooks, especially those that emphasize theory. Throughout, we have tried to include topics that a statistician would use in the practice of statistics and to cover these thoroughly. We don't develop every aspect of statistical theory; for example, we cover very little of the limit theorems in statistics (convergence in probability, convergence in distribution, almost sure convergence, Slutsky's theorem, etc.). We don't cover the Cramer–Rao lower bound or the Rao–Blackwell theorem. We cover joint continuous distributions using multiple integration, but we do not go into great depth.

The emphasis is on modern methods of statistical inference. We develop enough theory so that students will understand these methods. If a statistician or data scientist is to work effectively with practitioners, it is up to the statistician to be the one to explain how methods work, what assumptions underlie the methods, what the limitations are, and how (or whether) the assumptions can be checked. Subject matter experts (i.e., the nonstatisticians) are not trained to do this. This is why it is important for students of statistics to understand the underlying theory behind the methods.

The flip side of our approach is that we do not develop theory for theory's sake. No theory is developed for the purpose that it might be usable in a future course. We have found that students

who understand probability and the foundational concepts of statistical theory can understand and use advanced statistical methods. Without a solid grounding on the theory and concepts of statistics it is difficult to pick up new methods.

Calculus is used in a number of places in the book, so students will need at least one or two semesters of calculus. There are a few uses of multiple integrals when we discuss joint continuous distributions, and for these the third semester of calculus will be needed. An instructor can skip these topics or sidestep the use of multiple integrals. We use calculus when it is necessary, for example in getting expected values of continuous random variables. We use R throughout the book. Although we do cover an introduction to R, it would be helpful if students had some prior background in R.

We use data extensively throughout the book. Most of the data sets are real (although at times we give small data sets to introduce a method). Many of these data sets are large. In most cases, we have provided a csv (comma separated values) file for the data. We also provide the R code used in the book to analyze the data sets that we provide. This can be found at: www.cambridge.org/ProbStatsforDS

While the book's website contains information about getting R up and running, we offer the following advice about loading in data sets and packages. First, it is always good practice to set the working directory to the directory on your computer that contains your data files. You can do this with the `setwd()` command. For example,

```
setwd("C:/Users/Documents/Rfiles")
```

will force R to read (write) files from (to) this directory. Note two things: (1) the path must be enclosed in quotes, and (2) subdirectories are indicated by forward slashes, not backslashes. Second, many of the methods we apply in this book require special R packages to run. These packages are collections of functions, dataframes, etc. Before you can use a package you must (1) install it, and (2) load it in during each R session. To install a package, such as `dplyr`, type

```
install.packages("dplyr")
```

Then, every time you start a new R session, you will have to load this package using

```
library(dplyr)
```

You need only install a package once on your computer, but you must call `library()` each time you begin an R session.

If you type `library` for a package you haven't installed, you will get an error. For example, if you haven't installed the `testassay` package and if you type `library(testassay)`, then you will get an error like this:

```
Error in library(testassay) : there is no package called 'testassay'
```

The remedy is to first install the package by typing `install.packages("testassay")` and then typing `library(testassay)`. If you ever get an error like the following

```
Error in arrange(df, y) : could not find function "arrange"
```

there is a good chance you forgot to load the package that contains the function `arrange()`, which is in the `dplyr` package. The remedy is to first type `library(dplyr)`.

Most two-semester courses will include a fairly standard first semester, which would likely cover the following chapters:

Semester 1

Chapter	Topic
1	Introduction
2	Data Visualization
3	Basic Probability
4	Random Variables
5	Discrete Distributions
6	Continuous Distributions
7	About Data and Data Collection
8	Sampling Distributions
9	Point Estimation
10	Confidence Intervals
11	Hypothesis Testing
12	Hypothesis Tests for Two or More Populations

The choice of topics for a second course would depend on the nature of the course. For example, our book could be used in a mathematical statistics course that emphasizes applications of statistics without sacrificing any of the underlying theory. Such a course could use the following material in the second term:

Semester 2

Chapter	Topic
13	Hypothesis Tests for Categorical Data
14	Regression
15	Bayesian Methods
17	The Jackknife and Bootstrap
18	Generalized Linear Models and Regression Trees
20	Large-Scale Hypothesis Testing

For a course that leans toward data science, the second semester coverage might include:

Semester 2

Chapter	Topic
13	Hypothesis Tests for Categorical Data
14	Regression
16	Time Series Methods
17	The Jackknife and Bootstrap
18	Generalized Linear Models and Regression Trees
19	Cross-Validation and Estimates of Prediction Error
20	Large-Scale Hypothesis Testing

A course for scientists or engineers could include selected topics in the above chapters, with additional methods from Chapter 15. For example, a course in biostatistics might emphasize the sections on logistic regression, discrimination, and classification since these are frequently used in medical and public health research. Such a course could minimize or skip material on regression trees. Instructors

could also use this as a textbook for a one-semester course by selecting (and omitting) material in the early part of the book. For example, the following chapters could be covered in a one-semester course:

One-semester course emphasizing statistics

Chapter	Topic
1	Introduction
2	Data Visualization (omitting data visualization for survey data, geospatial data, and network data)
3	Basic Probability
4	Random Variables
5	Discrete Distributions (possibly omitting the hypergeometric and multinomial distributions)
6	Continuous Distributions (possibly skipping the Weibull, Beta distributions, and the sections on transformations, moment-generating functions, and QQ plots)
7	About Data and Data Collection (hitting just the main ideas)
8	Sampling Distributions (skipping the proof of the Central Limit Theorem)
9	Point Estimation
10	Confidence Intervals
11	Hypothesis Testing
12	Hypothesis Tests for Two or More Populations
13 or 14	Hypothesis Tests for Categorical Data/Regression

For situations where students have had a prior course on statistics (possibly one that did not use calculus), a course could be designed to emphasize data science:

One-semester course emphasizing data science

Chapter	Topics
4–6	Select topics in these chapters to bring students up to speed
7	About Data and Data Collection (hitting just the main ideas)
8	Sampling Distributions (skipping the proof of the Central Limit Theorem)
9	Point Estimation
10	Confidence Intervals
11	Hypothesis Testing
12	Hypothesis Tests for Two or More Populations
13	Hypothesis Tests for Categorical Data
14	Regression
17.	The Jackknife and Bootstrap
18.	Generalized Linear Models and Regression Trees
20.	Large-Scale Hypothesis Testing

This book was typeset in \LaTeX using a modified version of The Legrand Orange Book template originally created by Mathias Legrand and modified by Vel and the authors.

We would like to thank Emily Rigdon for \LaTeX ing much of the material in the book and Gary Smith for his careful reading and editing of the manuscript. We would also like to thank the staff at Cambridge, especially Lauren Cowles, Maggie Jeffers, and Lucy Edwards for their help in molding this book into what it has become, and for their patience through the process.

Steven E. Rigdon

Ronald D. Fricker, Jr.

Douglas C. Montgomery