

GLOBAL  
EDITION



# Introduction to Econometrics

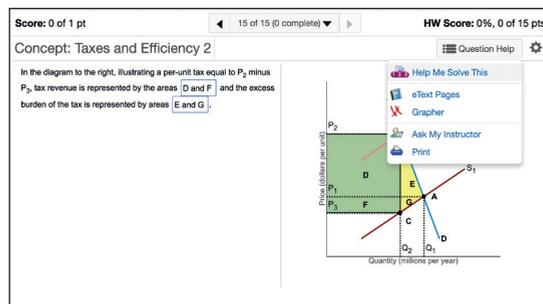
FOURTH EDITION

James H. Stock • Mark W. Watson



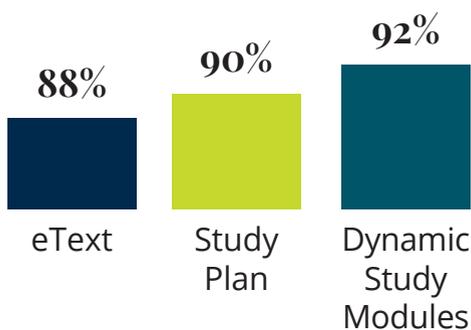
## Question Help

MyLab Economics homework and practice questions are correlated to the textbook, and many generate algorithmically to give students unlimited opportunity for mastery of concepts. If students get stuck, Learning Aids including Help Me Solve This and eText Pages walk them through the problem and identify helpful information in the text, giving them assistance when they need it most.



*"[MyLab Economics] provides ample practice and explanation of the concepts at hand."*

— Heather Burkett, University of Nebraska at Omaha



**% of students who found learning tool helpful**

**Dynamic Study Modules** help students study chapter topics effectively on their own by continuously assessing their **knowledge application** and performance in real time. These are available as prebuilt Prepare assignments, and are accessible on smartphones, tablets, and computers.

**Pearson eText** enhances student learning—both in and outside the classroom. Worked examples, videos, and interactive tutorials bring learning to life, while algorithmic practice and self-assessment opportunities test students' understanding of the material. Accessible anytime, anywhere via MyLab or the app.

The **MyLab Gradebook** offers an easy way for students and instructors to view course performance. Item Analysis allows instructors to quickly see trends by analyzing details like the number of students who answered correctly/incorrectly, time on task, and median time spend on a question by question basis. And because it's correlated with the AACSB Standards, instructors can track students' progress toward outcomes that the organization has deemed important in preparing students to be **leaders**.



**of students would tell their instructor to keep using MyLab Economics**

For additional details visit: [www.pearson.com/mylab/economics](http://www.pearson.com/mylab/economics)

# The Pearson Series in Economics

<b>Abel/Bernanke/Croushore</b> <i>Macroeconomics*</i> <sup>†</sup>	<b>Greene</b> <i>Econometric Analysis</i> <sup>†</sup>	<i>The Economics of Money, Banking, and Financial Markets, Business School Edition*</i>
<b>Acemoglu/Laibson/List</b> <i>Economics*</i> <sup>†</sup>	<b>Gregory/Stuart</b> <i>Russian and Soviet Economic Performance and Structure</i>	<i>Macroeconomics: Policy and Practice*</i>
<b>Bade/Parkin</b> <i>Foundations of Economics*</i> <sup>†</sup>	<b>Hartwick/Olewiler</b> <i>The Economics of Natural Resource Use</i>	<b>Murray</b> <i>Econometrics: A Modern Introduction</i>
<b>Berck/Helfand</b> <i>The Economics of the Environment</i>	<b>Heilbroner/Milberg</b> <i>The Making of the Economic Society</i>	<b>O'Sullivan/Sheffrin/Perez</b> <i>Economics: Principles, Applications and Tools*</i> <sup>†</sup>
<b>Bierman/Fernandez</b> <i>Game Theory with Economic Applications</i>	<b>Heyne/Boettke/Prychitko</b> <i>The Economic Way of Thinking</i>	<b>Parkin</b> <i>Economics*</i> <sup>†</sup>
<b>Blair/Rush</b> <i>The Economics of Managerial Decisions*</i> <sup>†</sup>	<b>Hubbard/O'Brien</b> <i>Economics*</i> <sup>†</sup>	<b>Perloff</b> <i>Microeconomics*</i> <sup>†</sup>
<b>Blanchard</b> <i>Macroeconomics*</i> <sup>†</sup>	<i>InEcon</i>	<i>Microeconomics: Theory and Applications with Calculus*</i> <sup>†</sup>
<b>Boyer</b> <i>Principles of Transportation Economics</i>	<i>Money, Banking, and the Financial System*</i>	<b>Pindyck/Rubinfeld</b> <i>Microeconomics*</i> <sup>†</sup>
<b>Brander/Perloff</b> <i>Managerial Economics and Strategy*</i> <sup>†</sup>	<b>Hubbard/O'Brien/Rafferty</b> <i>Macroeconomics*</i>	<b>Riddell/Shackelford/Stamos/Schneider</b> <i>Economics: A Tool for Critically Understanding Society</i>
<b>Branson</b> <i>Macroeconomic Theory and Policy</i>	<b>Hughes/Cain</b> <i>American Economic History</i>	<b>Roberts</b> <i>The Choice: A Fable of Free Trade and Protection</i>
<b>Bruce</b> <i>Public Finance and the American Economy</i>	<b>Husted/Melvin</b> <i>International Economics</i>	<b>Scherer</b> <i>Industry Structure, Strategy, and Public Policy</i>
<b>Carlton/Perloff</b> <i>Modern Industrial Organization</i>	<b>Jehle/Reny</b> <i>Advanced Microeconomic Theory</i>	<b>Schiller</b> <i>The Economics of Poverty and Discrimination</i>
<b>Case/Fair/Oster</b> <i>Principles of Economics*</i> <sup>†</sup>	<b>Keat/Young/Erflle</b> <i>Managerial Economics</i>	<b>Sherman</b> <i>Market Regulation</i>
<b>Chapman</b> <i>Environmental Economics: Theory, Application, and Policy</i>	<b>Klein</b> <i>Mathematical Methods for Economics</i>	<b>Stock/Watson</b> <i>Introduction to Econometrics</i> <sup>†</sup>
<b>Daniels/VanHoose</b> <i>International Monetary &amp; Financial Economics</i>	<b>Krugman/Obstfeld/Melitz</b> <i>International Economics: Theory &amp; Policy*</i> <sup>†</sup>	<b>Studenmund</b> <i>A Practical Guide to Using Econometrics</i> <sup>†</sup>
<b>Downs</b> <i>An Economic Theory of Democracy</i>	<b>Laidler</b> <i>The Demand for Money</i>	<b>Todaro/Smith</b> <i>Economic Development</i>
<b>Farnham</b> <i>Economics for Managers</i>	<b>Lynn</b> <i>Economic Development: Theory and Practice for a Divided World</i>	<b>Walters/Walters/Appel/Callahan/Centanni/Maex/O'Neill</b> <i>Conversations: Today's Students Discuss Today's Issues</i>
<b>Froyen</b> <i>Macroeconomics: Theories and Policies</i>	<b>Miller</b> <i>Economics Today*</i>	<b>Williamson</b> <i>Macroeconomics</i> <sup>†</sup>
<b>Fusfeld</b> <i>The Age of the Economist</i>	<b>Miller/Benjamin</b> <i>The Economics of Macro Issues</i>	
<b>Gerber</b> <i>International Economics*</i> <sup>†</sup>	<b>Miller/Benjamin/North</b> <i>The Economics of Public Issues</i>	
<b>Gordon</b> <i>Macroeconomics*</i>	<b>Mishkin</b> <i>The Economics of Money, Banking, and Financial Markets*</i> <sup>†</sup>	

\*denotes **MyLab Economics** titles. Visit [www.pearson.com/mylab/economics](http://www.pearson.com/mylab/economics) to learn more.

<sup>†</sup>denotes **Global Edition** titles.

# Introduction to Econometrics

FOURTH EDITION

GLOBAL EDITION

**James H. Stock**

Harvard University

**Mark W. Watson**

Princeton University



**Pearson**

---

Harlow, England • London • New York • Boston • San Francisco • Toronto • Sydney • Dubai • Singapore • Hong Kong  
Tokyo • Seoul • Taipei • New Delhi • Cape Town • Sao Paulo • Mexico City • Madrid • Amsterdam • Munich • Paris • Milan

Vice President, Business, Economics, and UK Courseware: Donna Battista  
Director of Portfolio Management: Adrienne D'Ambrosio  
Specialist Portfolio Manager: David Alexander  
Editorial Assistant: Nicole Nedwidek  
Project Editor, Global Edition: Paromita Banerjee  
Project Editor, Global Edition: Punita Kaur Mann  
Vice President, Product Marketing: Roxanne McCarley  
Product Marketing Assistant: Marianela Silvestri  
Manager of Field Marketing, Business Publishing: Adam Goldstein  
Executive Field Marketing Manager: Carlie Marvel  
Vice President, Production and Digital Studio, Arts and Business: Etain O'Dea  
Director, Production and Digital Studio, Business and Economics: Ashley Santora  
Managing Producer, Business: Alison Kalil  
Content Producer: Christine Donovan  
Content Producer, Global Edition: Nikhil Rakshit

Operations Specialist: Carol Melville  
Senior Manufacturing Controller, Global Edition: Kay Holman  
Manager, Learning Tools: Brian Surette  
Senior Learning Tools Strategist: Emily Biberger  
Managing Producer, Digital Studio and GLP: James Bateman  
Managing Producer, Digital Studio: Diane Lombardo  
Digital Studio Producer: Melissa Honig  
Digital Studio Producer: Alana Coles  
Digital Content Team Lead: Noel Lotz  
Digital Content Project Lead: Noel Lotz  
Manager, Media Production, Global Edition: Vikram Kumar  
Project Manager: Vikash Sharma, Cenveo Publisher Services  
Interior Design: Cenveo Publisher Services  
Cover Design: Lumina Datamatics  
Cover Art: GarryKillian / Shutterstock

Acknowledgments of third-party content appear on the appropriate page within the text.

Pearson Education Limited  
KAO Two  
KAO Park  
Harlow  
CM17 9NA  
United Kingdom

and Associated Companies throughout the world

Visit us on the World Wide Web at: [www.pearsonglobaleditions.com](http://www.pearsonglobaleditions.com)

© Pearson Education Limited 2020

The rights of James H. Stock and Mark W. Watson, to be identified as the authors of this work, have been asserted by them in accordance with the Copyright, Designs and Patents Act 1988.

Authorized adaptation from the United States edition, entitled *Introduction to Econometrics*, 4th Edition, ISBN 978-0-13-446199-1 by James H. Stock and Mark W. Watson, published by Pearson Education © 2020.

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, without either the prior written permission of the publisher or a license permitting restricted copying in the United Kingdom issued by the Copyright Licensing Agency Ltd, Saffron House, 6–10 Kirby Street, London EC1N 8TS.

All trademarks used herein are the property of their respective owners. The use of any trademark in this text does not vest in the author or publisher any trademark ownership rights in such trademarks, nor does the use of such trademarks imply any affiliation with or endorsement of this book by such owners. For information regarding permissions, request forms, and the appropriate contacts within the Pearson Education Global Rights and Permissions department, please visit [www.pearsoned.com/permissions/](http://www.pearsoned.com/permissions/).

This eBook is a standalone product and may or may not include all assets that were part of the print version. It also does not provide access to other Pearson digital products like MyLab and Mastering. The publisher reserves the right to remove any material in this eBook at any time.

#### **British Library Cataloguing-in-Publication Data**

A catalogue record for this book is available from the British Library

**ISBN 10:** 1-292-26445-4

**ISBN 13:** 978-1-292-26445-5

**eBook ISBN 13:** 978-1-292-26452-3

Typeset in Times NR MT Pro by Cenveo® Publisher Services

# Brief Contents

## **PART ONE Introduction and Review**

---

Chapter 1	Economic Questions and Data	43
Chapter 2	Review of Probability	55
Chapter 3	Review of Statistics	103

## **PART TWO Fundamentals of Regression Analysis**

---

Chapter 4	Linear Regression with One Regressor	143
Chapter 5	Regression with a Single Regressor: Hypothesis Tests and Confidence Intervals	178
Chapter 6	Linear Regression with Multiple Regressors	211
Chapter 7	Hypothesis Tests and Confidence Intervals in Multiple Regression	247
Chapter 8	Nonlinear Regression Functions	277
Chapter 9	Assessing Studies Based on Multiple Regression	330

## **PART THREE Further Topics in Regression Analysis**

---

Chapter 10	Regression with Panel Data	361
Chapter 11	Regression with a Binary Dependent Variable	392
Chapter 12	Instrumental Variables Regression	427
Chapter 13	Experiments and Quasi-Experiments	474
Chapter 14	Prediction with Many Regressors and Big Data	514

## **PART FOUR Regression Analysis of Economic Time Series Data**

---

Chapter 15	Introduction to Time Series Regression and Forecasting	554
Chapter 16	Estimation of Dynamic Causal Effects	609
Chapter 17	Additional Topics in Time Series Regression	649

## **PART FIVE Regression Analysis of Economic Time Series Data**

---

Chapter 18	The Theory of Linear Regression with One Regressor	687
Chapter 19	The Theory of Multiple Regression	713

This page intentionally left blank

# Contents

*Preface* 27

## **PART ONE Introduction and Review**

---

### **CHAPTER 1 Economic Questions and Data 43**

- 1.1 Economic Questions We Examine 43
  - Question #1: Does Reducing Class Size Improve Elementary School Education? 43
  - Question #2: Is There Racial Discrimination in the Market for Home Loans? 44
  - Question #3: Does Healthcare Spending Improve Health Outcomes? 45
  - Question #4: By How Much Will U.S. GDP Grow Next Year? 46
  - Quantitative Questions, Quantitative Answers 47
- 1.2 Causal Effects and Idealized Experiments 47
  - Estimation of Causal Effects 48
  - Prediction, Forecasting, and Causality 48
- 1.3 Data: Sources and Types 49
  - Experimental versus Observational Data 49
  - Cross-Sectional Data 50
  - Time Series Data 51
  - Panel Data 52

### **CHAPTER 2 Review of Probability 55**

- 2.1 Random Variables and Probability Distributions 56
  - Probabilities, the Sample Space, and Random Variables 56
  - Probability Distribution of a Discrete Random Variable 56
  - Probability Distribution of a Continuous Random Variable 58
- 2.2 Expected Values, Mean, and Variance 60
  - The Expected Value of a Random Variable 60
  - The Standard Deviation and Variance 61
  - Mean and Variance of a Linear Function of a Random Variable 62
  - Other Measures of the Shape of a Distribution 63
  - Standardized Random Variables 65
- 2.3 Two Random Variables 65
  - Joint and Marginal Distributions 65
  - Conditional Distributions 66
  - Independence 70
  - Covariance and Correlation 70
  - The Mean and Variance of Sums of Random Variables 71

2.4	The Normal, Chi-Squared, Student $t$ , and $F$ Distributions	75
	The Normal Distribution	75
	The Chi-Squared Distribution	80
	The Student $t$ Distribution	80
	The $F$ Distribution	80
2.5	Random Sampling and the Distribution of the Sample Average	81
	Random Sampling	81
	The Sampling Distribution of the Sample Average	82
2.6	Large-Sample Approximations to Sampling Distributions	85
	The Law of Large Numbers and Consistency	85
	The Central Limit Theorem	86
	APPENDIX 2.1 Derivation of Results in Key Concept 2.3	100
	APPENDIX 2.2 The Conditional Mean as the Minimum Mean Squared Error Predictor	101

### CHAPTER 3 Review of Statistics 103

3.1	Estimation of the Population Mean	104
	Estimators and Their Properties	104
	Properties of $\bar{Y}$	106
	The Importance of Random Sampling	108
3.2	Hypothesis Tests Concerning the Population Mean	109
	Null and Alternative Hypotheses	109
	The $p$ -Value	110
	Calculating the $p$ -Value When $\sigma_Y$ Is Known	111
	The Sample Variance, Sample Standard Deviation, and Standard Error	112
	Calculating the $p$ -Value When $\sigma_Y$ Is Unknown	113
	The $t$ -Statistic	113
	Hypothesis Testing with a Prespecified Significance Level	114
	One-Sided Alternatives	116
3.3	Confidence Intervals for the Population Mean	117
3.4	Comparing Means from Different Populations	119
	Hypothesis Tests for the Difference Between Two Means	119
	Confidence Intervals for the Difference Between Two Population Means	120
3.5	Differences-of-Means Estimation of Causal Effects Using Experimental Data	121
	The Causal Effect as a Difference of Conditional Expectations	121
	Estimation of the Causal Effect Using Differences of Means	121
3.6	Using the $t$ -Statistic When the Sample Size Is Small	123
	The $t$ -Statistic and the Student $t$ Distribution	125
	Use of the Student $t$ Distribution in Practice	126

- 3.7 Scatterplots, the Sample Covariance, and the Sample Correlation 127
  - Scatterplots 127
  - Sample Covariance and Correlation 127
  - APPENDIX 3.1 The U.S. Current Population Survey 141
  - APPENDIX 3.2 Two Proofs That  $\bar{Y}$  Is the Least Squares Estimator of  $\mu_Y$  141
  - APPENDIX 3.3 A Proof That the Sample Variance Is Consistent 142

## PART TWO Fundamentals of Regression Analysis

---

### CHAPTER 4 Linear Regression with One Regressor 143

- 4.1 The Linear Regression Model 144
- 4.2 Estimating the Coefficients of the Linear Regression Model 147
  - The Ordinary Least Squares Estimator 148
  - OLS Estimates of the Relationship Between Test Scores and the Student-Teacher Ratio 149
  - Why Use the OLS Estimator? 151
- 4.3 Measures of Fit and Prediction Accuracy 153
  - The  $R^2$  153
  - The Standard Error of the Regression 154
  - Prediction Using OLS 155
  - Application to the Test Score Data 155
- 4.4 The Least Squares Assumptions for Causal Inference 156
  - Assumption 1: The Conditional Distribution of  $u_i$  Given  $X_i$  Has a Mean of Zero 157
  - Assumption 2:  $(X_i, Y_i), i = 1, \dots, n$ , Are Independently and Identically Distributed 158
  - Assumption 3: Large Outliers Are Unlikely 159
  - Use of the Least Squares Assumptions 160
- 4.5 The Sampling Distribution of the OLS Estimators 161
- 4.6 Conclusion 164
  - APPENDIX 4.1 The California Test Score Data Set 172
  - APPENDIX 4.2 Derivation of the OLS Estimators 172
  - APPENDIX 4.3 Sampling Distribution of the OLS Estimator 173
  - APPENDIX 4.4 The Least Squares Assumptions for Prediction 176

### CHAPTER 5 Regression with a Single Regressor: Hypothesis Tests and Confidence Intervals 178

- 5.1 Testing Hypotheses About One of the Regression Coefficients 178
  - Two-Sided Hypotheses Concerning  $\beta_1$  179
  - One-Sided Hypotheses Concerning  $\beta_1$  182
  - Testing Hypotheses About the Intercept  $\beta_0$  184
- 5.2 Confidence Intervals for a Regression Coefficient 184

- 5.3 Regression When  $X$  Is a Binary Variable 186
  - Interpretation of the Regression Coefficients 186
- 5.4 Heteroskedasticity and Homoskedasticity 188
  - What Are Heteroskedasticity and Homoskedasticity? 188
  - Mathematical Implications of Homoskedasticity 190
  - What Does This Mean in Practice? 192
- 5.5 The Theoretical Foundations of Ordinary Least Squares 194
  - Linear Conditionally Unbiased Estimators and the Gauss–Markov Theorem 194
  - Regression Estimators Other Than OLS 195
- 5.6 Using the  $t$ -Statistic in Regression When the Sample Size Is Small 196
  - The  $t$ -Statistic and the Student  $t$  Distribution 196
  - Use of the Student  $t$  Distribution in Practice 197
- 5.7 Conclusion 197
  - APPENDIX 5.1 Formulas for OLS Standard Errors 206
  - APPENDIX 5.2 The Gauss–Markov Conditions and a Proof of the Gauss–Markov Theorem 207

## CHAPTER 6 Linear Regression with Multiple Regressors 211

- 6.1 Omitted Variable Bias 211
  - Definition of Omitted Variable Bias 212
  - A Formula for Omitted Variable Bias 214
  - Addressing Omitted Variable Bias by Dividing the Data into Groups 215
- 6.2 The Multiple Regression Model 217
  - The Population Regression Line 217
  - The Population Multiple Regression Model 218
- 6.3 The OLS Estimator in Multiple Regression 220
  - The OLS Estimator 220
  - Application to Test Scores and the Student–Teacher Ratio 221
- 6.4 Measures of Fit in Multiple Regression 222
  - The Standard Error of the Regression ( $SER$ ) 222
  - The  $R^2$  223
  - The Adjusted  $R^2$  223
  - Application to Test Scores 224
- 6.5 The Least Squares Assumptions for Causal Inference in Multiple Regression 225
  - Assumption 1: The Conditional Distribution of  $u_i$  Given  $X_{1i}, X_{2i}, \dots, X_{ki}$  Has a Mean of 0 225
  - Assumption 2:  $(X_{1i}, X_{2i}, \dots, X_{ki}, Y_i), i = 1, \dots, n$ , Are i.i.d. 225
  - Assumption 3: Large Outliers Are Unlikely 225
  - Assumption 4: No Perfect Multicollinearity 226

- 6.6 The Distribution of the OLS Estimators in Multiple Regression 227
- 6.7 Multicollinearity 228
  - Examples of Perfect Multicollinearity 228
  - Imperfect Multicollinearity 230
- 6.8 Control Variables and Conditional Mean Independence 231
  - Control Variables and Conditional Mean Independence 232
- 6.9 Conclusion 234
  - APPENDIX 6.1 Derivation of Equation (6.1) 242
  - APPENDIX 6.2 Distribution of the OLS Estimators When There Are Two Regressors and Homoskedastic Errors 243
  - APPENDIX 6.3 The Frisch–Waugh Theorem 243
  - APPENDIX 6.4 The Least Squares Assumptions for Prediction with Multiple Regressors 244
  - APPENDIX 6.5 Distribution of OLS Estimators in Multiple Regression with Control Variables 245

## CHAPTER 7 Hypothesis Tests and Confidence Intervals in Multiple Regression 247

- 7.1 Hypothesis Tests and Confidence Intervals for a Single Coefficient 247
  - Standard Errors for the OLS Estimators 247
  - Hypothesis Tests for a Single Coefficient 248
  - Confidence Intervals for a Single Coefficient 249
  - Application to Test Scores and the Student–Teacher Ratio 249
- 7.2 Tests of Joint Hypotheses 251
  - Testing Hypotheses on Two or More Coefficients 252
  - The  $F$ -Statistic 253
  - Application to Test Scores and the Student–Teacher Ratio 255
  - The Homoskedasticity-Only  $F$ -Statistic 256
- 7.3 Testing Single Restrictions Involving Multiple Coefficients 258
- 7.4 Confidence Sets for Multiple Coefficients 259
- 7.5 Model Specification for Multiple Regression 260
  - Model Specification and Choosing Control Variables 261
  - Interpreting the  $R^2$  and the Adjusted  $R^2$  in Practice 262
- 7.6 Analysis of the Test Score Data Set 262
- 7.7 Conclusion 268
  - APPENDIX 7.1 The Bonferroni Test of a Joint Hypothesis 274

<b>CHAPTER 8</b>	<b>Nonlinear Regression Functions</b>	<b>277</b>
8.1	A General Strategy for Modeling Nonlinear Regression Functions	279
	Test Scores and District Income	279
	The Effect on Y of a Change in X in Nonlinear Specifications	282
	A General Approach to Modeling Nonlinearities Using Multiple Regression	285
8.2	Nonlinear Functions of a Single Independent Variable	286
	Polynomials	286
	Logarithms	288
	Polynomial and Logarithmic Models of Test Scores and District Income	296
8.3	Interactions Between Independent Variables	297
	Interactions Between Two Binary Variables	298
	Interactions Between a Continuous and a Binary Variable	300
	Interactions Between Two Continuous Variables	305
8.4	Nonlinear Effects on Test Scores of the Student–Teacher Ratio	310
	Discussion of Regression Results	310
	Summary of Findings	314
8.5	Conclusion	315
	<a href="#">APPENDIX 8.1</a> Regression Functions That Are Nonlinear in the Parameters	325
	<a href="#">APPENDIX 8.2</a> Slopes and Elasticities for Nonlinear Regression Functions	328
<b>CHAPTER 9</b>	<b>Assessing Studies Based on Multiple Regression</b>	<b>330</b>
9.1	Internal and External Validity	330
	Threats to Internal Validity	331
	Threats to External Validity	332
9.2	Threats to Internal Validity of Multiple Regression Analysis	333
	Omitted Variable Bias	334
	Misspecification of the Functional Form of the Regression Function	336
	Measurement Error and Errors-in-Variables Bias	336
	Missing Data and Sample Selection	339
	Simultaneous Causality	341
	Sources of Inconsistency of OLS Standard Errors	343
9.3	Internal and External Validity When the Regression Is Used for Prediction	344
9.4	Example: Test Scores and Class Size	345
	External Validity	346
	Internal Validity	352
	Discussion and Implications	353
9.5	Conclusion	354
	<a href="#">APPENDIX 9.1</a> The Massachusetts Elementary School Testing Data	360

## **PART THREE Further Topics in Regression Analysis**

---

### **CHAPTER 10 Regression with Panel Data 361**

- 10.1 Panel Data 362
  - Example: Traffic Deaths and Alcohol Taxes 362
- 10.2 Panel Data with Two Time Periods: “Before and After” Comparisons 365
- 10.3 Fixed Effects Regression 367
  - The Fixed Effects Regression Model 367
  - Estimation and Inference 369
  - Application to Traffic Deaths 370
- 10.4 Regression with Time Fixed Effects 371
  - Time Effects Only 371
  - Both Entity and Time Fixed Effects 372
- 10.5 The Fixed Effects Regression Assumptions and Standard Errors for Fixed Effects Regression 374
  - The Fixed Effects Regression Assumptions 374
  - Standard Errors for Fixed Effects Regression 376
- 10.6 Drunk Driving Laws and Traffic Deaths 377
- 10.7 Conclusion 381
  - [APPENDIX 10.1 The State Traffic Fatality Data Set](#) 387
  - [APPENDIX 10.2 Standard Errors for Fixed Effects Regression](#) 388

### **CHAPTER 11 Regression with a Binary Dependent Variable 392**

- 11.1 Binary Dependent Variables and the Linear Probability Model 393
  - Binary Dependent Variables 393
  - The Linear Probability Model 395
- 11.2 Probit and Logit Regression 397
  - Probit Regression 397
  - Logit Regression 401
  - Comparing the Linear Probability, Probit, and Logit Models 403
- 11.3 Estimation and Inference in the Logit and Probit Models 404
  - Nonlinear Least Squares Estimation 404
  - Maximum Likelihood Estimation 405
  - Measures of Fit 406
- 11.4 Application to the Boston HMDA Data 407
- 11.5 Conclusion 413
  - [APPENDIX 11.1 The Boston HMDA Data Set](#) 421
  - [APPENDIX 11.2 Maximum Likelihood Estimation](#) 421
  - [APPENDIX 11.3 Other Limited Dependent Variable Models](#) 424

<b>CHAPTER 12</b>	<b>Instrumental Variables Regression</b>	<b>427</b>
12.1	The IV Estimator with a Single Regressor and a Single Instrument	428
	The IV Model and Assumptions	428
	The Two Stage Least Squares Estimator	429
	Why Does IV Regression Work?	429
	The Sampling Distribution of the TSLS Estimator	434
	Application to the Demand for Cigarettes	435
12.2	The General IV Regression Model	437
	TSLS in the General IV Model	439
	Instrument Relevance and Exogeneity in the General IV Model	440
	The IV Regression Assumptions and Sampling Distribution of the TSLS Estimator	441
	Inference Using the TSLS Estimator	442
	Application to the Demand for Cigarettes	443
12.3	Checking Instrument Validity	444
	Assumption 1: Instrument Relevance	444
	Assumption 2: Instrument Exogeneity	446
12.4	Application to the Demand for Cigarettes	450
12.5	Where Do Valid Instruments Come From?	454
	Three Examples	455
12.6	Conclusion	459
	<b>APPENDIX 12.1</b> The Cigarette Consumption Panel Data Set	467
	<b>APPENDIX 12.2</b> Derivation of the Formula for the TSLS Estimator in Equation (12.4)	467
	<b>APPENDIX 12.3</b> Large-Sample Distribution of the TSLS Estimator	468
	<b>APPENDIX 12.4</b> Large-Sample Distribution of the TSLS Estimator When the Instrument Is Not Valid	469
	<b>APPENDIX 12.5</b> Instrumental Variables Analysis with Weak Instruments	470
	<b>APPENDIX 12.6</b> TSLS with Control Variables	472
<b>CHAPTER 13</b>	<b>Experiments and Quasi-Experiments</b>	<b>474</b>
13.1	Potential Outcomes, Causal Effects, and Idealized Experiments	475
	Potential Outcomes and the Average Causal Effect	475
	Econometric Methods for Analyzing Experimental Data	476
13.2	Threats to Validity of Experiments	478
	Threats to Internal Validity	478
	Threats to External Validity	481
13.3	Experimental Estimates of the Effect of Class Size Reductions	482
	Experimental Design	482
	Analysis of the STAR Data	483
	Comparison of the Observational and Experimental Estimates of Class Size Effects	488

13.4	Quasi-Experiments	490
	Examples	490
	The Differences-in-Differences Estimator	492
	Instrumental Variables Estimators	494
	Regression Discontinuity Estimators	495
13.5	Potential Problems with Quasi-Experiments	496
	Threats to Internal Validity	496
	Threats to External Validity	498
13.6	Experimental and Quasi-Experimental Estimates in Heterogeneous Populations	498
	OLS with Heterogeneous Causal Effects	499
	IV Regression with Heterogeneous Causal Effects	500
13.7	Conclusion	503
	APPENDIX 13.1 The Project STAR Data Set	510
	APPENDIX 13.2 IV Estimation When the Causal Effect Varies Across Individuals	511
	APPENDIX 13.3 The Potential Outcomes Framework for Analyzing Data from Experiments	512
<b>CHAPTER 14</b>	<b>Prediction with Many Regressors and Big Data</b>	<b>514</b>
14.1	What Is “Big Data”?	515
14.2	The Many-Predictor Problem and OLS	516
	The Mean Squared Prediction Error	518
	The First Least Squares Assumption for Prediction	519
	The Predictive Regression Model with Standardized Regressors	519
	The MSPE of OLS and the Principle of Shrinkage	521
	Estimation of the MSPE	522
14.3	Ridge Regression	524
	Shrinkage via Penalization and Ridge Regression	524
	Estimation of the Ridge Shrinkage Parameter by Cross Validation	525
	Application to School Test Scores	526
14.4	The Lasso	527
	Shrinkage Using the Lasso	528
	Application to School Test Scores	531
14.5	Principal Components	532
	Principals Components with Two Variables	532
	Principal Components with $k$ Variables	534
	Application to School Test Scores	536
14.6	Predicting School Test Scores with Many Predictors	537

- 14.7 Conclusion 542
  - APPENDIX 14.1 The California School Test Score Data Set 551
  - APPENDIX 14.2 Derivation of Equation (14.4) for  $k = 1$  551
  - APPENDIX 14.3 The Ridge Regression Estimator When  $k = 1$  551
  - APPENDIX 14.4 The Lasso Estimator When  $k = 1$  552
  - APPENDIX 14.5 Computing Out-of-Sample Predictions in the Standardized Regression Model 552

## **PART FOUR Regression Analysis of Economic Time Series Data**

---

### **CHAPTER 15 Introduction to Time Series Regression and Forecasting 554**

- 15.1 Introduction to Time Series Data and Serial Correlation 555
  - Real GDP in the United States 555
  - Lags, First Differences, Logarithms, and Growth Rates 555
  - Autocorrelation 558
  - Other Examples of Economic Time Series 560
- 15.2 Stationarity and the Mean Squared Forecast Error 561
  - Stationarity 561
  - Forecasts and Forecast Errors 562
  - The Mean Squared Forecast Error 563
- 15.3 Autoregressions 565
  - The First-Order Autoregressive Model 565
  - The  $p^{\text{th}}$ -Order Autoregressive Model 567
- 15.4 Time Series Regression with Additional Predictors and the Autoregressive Distributed Lag Model 568
  - Forecasting GDP Growth Using the Term Spread 569
  - The Autoregressive Distributed Lag Model 570
  - The Least Squares Assumptions for Forecasting with Multiple Predictors 571
- 15.5 Estimation of the MSFE and Forecast Intervals 573
  - Estimation of the MSFE 573
  - Forecast Uncertainty and Forecast Intervals 576
- 15.6 Estimating the Lag Length Using Information Criteria 578
  - Determining the Order of an Autoregression 578
  - Lag Length Selection in Time Series Regression with Multiple Predictors 581
- 15.7 Nonstationarity I: Trends 582
  - What Is a Trend? 582
  - Problems Caused by Stochastic Trends 584
  - Detecting Stochastic Trends: Testing for a Unit AR Root 586
  - Avoiding the Problems Caused by Stochastic Trends 588

15.8	Nonstationarity II: Breaks	589
	What Is a Break?	589
	Testing for Breaks	589
	Detecting Breaks Using Pseudo Out-of-Sample Forecasts	594
	Avoiding the Problems Caused by Breaks	595
15.9	Conclusion	596
	APPENDIX 15.1 Time Series Data Used in Chapter 15	604
	APPENDIX 15.2 Stationarity in the AR(1) Model	605
	APPENDIX 15.3 Lag Operator Notation	606
	APPENDIX 15.4 ARMA Models	607
	APPENDIX 15.5 Consistency of the BIC Lag Length Estimator	607
<b>CHAPTER 16</b>	<b>Estimation of Dynamic Causal Effects</b>	<b>609</b>
16.1	An Initial Taste of the Orange Juice Data	610
16.2	Dynamic Causal Effects	612
	Causal Effects and Time Series Data	612
	Two Types of Exogeneity	615
16.3	Estimation of Dynamic Causal Effects with Exogenous Regressors	617
	The Distributed Lag Model Assumptions	617
	Autocorrelated $u_t$ , Standard Errors, and Inference	618
	Dynamic Multipliers and Cumulative Dynamic Multipliers	618
16.4	Heteroskedasticity- and Autocorrelation-Consistent Standard Errors	620
	Distribution of the OLS Estimator with Autocorrelated Errors	620
	HAC Standard Errors	621
16.5	Estimation of Dynamic Causal Effects with Strictly Exogenous Regressors	624
	The Distributed Lag Model with AR(1) Errors	625
	OLS Estimation of the ADL Model	627
	GLS Estimation	628
16.6	Orange Juice Prices and Cold Weather	630
16.7	Is Exogeneity Plausible? Some Examples	637
	U.S. Income and Australian Exports	637
	Oil Prices and Inflation	637
	Monetary Policy and Inflation	638
	The Growth Rate of GDP and the Term Spread	638
16.8	Conclusion	639
	APPENDIX 16.1 The Orange Juice Data Set	646
	APPENDIX 16.2 The ADL Model and Generalized Least Squares in Lag Operator Notation	647

**CHAPTER 17 Additional Topics in Time Series Regression 649**

- 17.1 Vector Autoregressions 649
  - The VAR Model 650
  - A VAR Model of the Growth Rate of GDP and the Term Spread 653
- 17.2 Multi-period Forecasts 654
  - Iterated Multi-period Forecasts 654
  - Direct Multi-period Forecasts 656
  - Which Method Should You Use? 658
- 17.3 Orders of Integration and the Nonnormality of Unit Root Test Statistics 658
  - Other Models of Trends and Orders of Integration 659
  - Why Do Unit Root Tests Have Nonnormal Distributions? 661
- 17.4 Cointegration 663
  - Cointegration and Error Correction 663
  - How Can You Tell Whether Two Variables Are Cointegrated? 664
  - Estimation of Cointegrating Coefficients 665
  - Extension to Multiple Cointegrated Variables 666
- 17.5 Volatility Clustering and Autoregressive Conditional Heteroskedasticity 667
  - Volatility Clustering 667
  - Realized Volatility 668
  - Autoregressive Conditional Heteroskedasticity 669
  - Application to Stock Price Volatility 670
- 17.6 Forecasting with Many Predictors Using Dynamic Factor Models and Principal Components 671
  - The Dynamic Factor Model 672
  - The DFM: Estimation and Forecasting 673
  - Application to U.S. Macroeconomic Data 676
- 17.7 Conclusion 682
  - [APPENDIX 17.1 The Quarterly U.S. Macro Data Set](#) 686

**PART FIVE Regression Analysis of Economic Time Series Data**

---

**CHAPTER 18 The Theory of Linear Regression with One Regressor 687**

- 18.1 The Extended Least Squares Assumptions and the OLS Estimator 688
  - The Extended Least Squares Assumptions 688
  - The OLS Estimator 689
- 18.2 Fundamentals of Asymptotic Distribution Theory 690
  - Convergence in Probability and the Law of Large Numbers 690
  - The Central Limit Theorem and Convergence in Distribution 692

	Slutsky's Theorem and the Continuous Mapping Theorem	693
	Application to the $t$ -Statistic Based on the Sample Mean	694
18.3	Asymptotic Distribution of the OLS Estimator and $t$ -Statistic	695
	Consistency and Asymptotic Normality of the OLS Estimators	695
	Consistency of Heteroskedasticity-Robust Standard Errors	695
	Asymptotic Normality of the Heteroskedasticity-Robust $t$ -Statistic	696
18.4	Exact Sampling Distributions When the Errors Are Normally Distributed	697
	Distribution of $\hat{\beta}_1$ with Normal Errors	697
	Distribution of the Homoskedasticity-Only $t$ -Statistic	698
18.5	Weighted Least Squares	699
	WLS with Known Heteroskedasticity	700
	WLS with Heteroskedasticity of Known Functional Form	701
	Heteroskedasticity-Robust Standard Errors or WLS?	703
	APPENDIX 18.1 The Normal and Related Distributions and Moments of Continuous Random Variables	709
	APPENDIX 18.2 Two Inequalities	711
<b>CHAPTER 19</b>	<b>The Theory of Multiple Regression</b>	<b>713</b>
19.1	The Linear Multiple Regression Model and OLS Estimator in Matrix Form	714
	The Multiple Regression Model in Matrix Notation	714
	The Extended Least Squares Assumptions	715
	The OLS Estimator	716
19.2	Asymptotic Distribution of the OLS Estimator and $t$ -Statistic	717
	The Multivariate Central Limit Theorem	718
	Asymptotic Normality of $\hat{\beta}$	718
	Heteroskedasticity-Robust Standard Errors	719
	Confidence Intervals for Predicted Effects	720
	Asymptotic Distribution of the $t$ -Statistic	720
19.3	Tests of Joint Hypotheses	721
	Joint Hypotheses in Matrix Notation	721
	Asymptotic Distribution of the $F$ -Statistic	721
	Confidence Sets for Multiple Coefficients	722
19.4	Distribution of Regression Statistics with Normal Errors	722
	Matrix Representations of OLS Regression Statistics	723
	Distribution of $\hat{\beta}$ with Independent Normal Errors	724
	Distribution of $s_e^2$	724
	Homoskedasticity-Only Standard Errors	724
	Distribution of the $t$ -Statistic	725
	Distribution of the $F$ -Statistic	725

19.5	Efficiency of the OLS Estimator with Homoskedastic Errors	726
	The Gauss–Markov Conditions for Multiple Regression	726
	Linear Conditionally Unbiased Estimators	726
	The Gauss–Markov Theorem for Multiple Regression	727
19.6	Generalized Least Squares	728
	The GLS Assumptions	729
	GLS When $\Omega$ Is Known	730
	GLS When $\Omega$ Contains Unknown Parameters	731
	The Conditional Mean Zero Assumption and GLS	731
19.7	Instrumental Variables and Generalized Method of Moments Estimation	733
	The IV Estimator in Matrix Form	733
	Asymptotic Distribution of the TSLS Estimator	734
	Properties of TSLS When the Errors Are Homoskedastic	735
	Generalized Method of Moments Estimation in Linear Models	738
	APPENDIX 19.1 Summary of Matrix Algebra	748
	APPENDIX 19.2 Multivariate Distributions	752
	APPENDIX 19.3 Derivation of the Asymptotic Distribution of $\hat{\beta}$	753
	APPENDIX 19.4 Derivations of Exact Distributions of OLS Test Statistics with Normal Errors	754
	APPENDIX 19.5 Proof of the Gauss–Markov Theorem for Multiple Regression	755
	APPENDIX 19.6 Proof of Selected Results for IV and GMM Estimation	756
	APPENDIX 19.7 Regression with Many Predictors: MSPE, Ridge Regression, and Principal Components Analysis	758
	<b>Appendix</b>	<b>763</b>
	<b>References</b>	<b>771</b>
	<b>Glossary</b>	<b>775</b>
	<b>Index</b>	<b>785</b>

# Key Concepts

## PART ONE Introduction and Review

---

- 1.1 Cross-Sectional, Time Series, and Panel Data 53
- 2.1 Expected Value and the Mean 60
- 2.2 Variance and Standard Deviation 61
- 2.3 Means, Variances, and Covariances of Sums of Random Variables 74
- 2.4 Computing Probabilities and Involving Normal Random Variables 76
- 2.5 Simple Random Sampling and i.i.d. Random Variables 82
- 2.6 Convergence in Probability, Consistency, and the Law of Large Numbers 86
- 2.7 The Central Limit Theorem 89
- 3.1 Estimators and Estimates 105
- 3.2 Bias, Consistency, and Efficiency 105
- 3.3 Efficiency of  $\bar{Y}$ :  $\bar{Y}$  Is BLUE 107
- 3.4 The Standard Error of  $\bar{Y}$  113
- 3.5 The Terminology of Hypothesis Testing 115
- 3.6 Testing the Hypothesis  $E(Y) = \mu_{Y,0}$  Against the Alternative  $E(Y) \neq \mu_{Y,0}$  116
- 3.7 Confidence Intervals for the Population Mean 118

## PART TWO Fundamentals of Regression Analysis

---

- 4.1 Terminology for the Linear Regression Model with a Single Regressor 146
- 4.2 The OLS Estimator, Predicted Values, and Residuals 150
- 4.3 The Least Squares Assumptions for Causal Inference 160
- 4.4 Large-Sample Distributions of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  162
- 5.1 General Form of the  $t$ -Statistic 179
- 5.2 Testing the Hypothesis  $\beta_1 = \beta_{1,0}$  Against the Alternative  $\beta_1 \neq \beta_{1,0}$  181
- 5.3 Confidence Interval for  $\beta_1$  185
- 5.4 Heteroskedasticity and Homoskedasticity 190
- 5.5 The Gauss–Markov Theorem for  $\hat{\beta}_1$  195
- 6.1 Omitted Variable Bias in Regression with a Single Regressor 213
- 6.2 The Multiple Regression Model 219
- 6.3 The OLS Estimators, Predicted Values, and Residuals in the Multiple Regression Model 221
- 6.4 The Least Squares Assumptions for Causal Inference in the Multiple Regression Model 227
- 6.5 Large-Sample Distribution of  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$  228
- 6.6 The Least Squares Assumptions for Causal Inference in the Multiple Regression Model with Control Variables 233
- 7.1 Testing the Hypothesis  $\beta_j = \beta_{j,0}$  Against the Alternative  $\beta_j \neq \beta_{j,0}$  249

- 7.2 Confidence Intervals for a Single Coefficient in Multiple Regression 250
- 7.3  $R^2$  and  $\bar{R}^2$ : What They Tell You—and What They Don't 263
- 8.1 The Expected Change in  $Y$  from a Change in  $X_1$  in the Nonlinear Regression Model [Equation (8.3)] 283
- 8.2 Logarithms in Regression: Three Cases 295
- 8.3 A Method for Interpreting Coefficients in Regressions with Binary Variables 299
- 8.4 Interactions Between Binary and Continuous Variables 302
- 8.5 Interactions in Multiple Regression 306
- 9.1 Internal and External Validity 331
- 9.2 Omitted Variable Bias: Should I Include More Variables in My Regression? 335
- 9.3 Functional Form Misspecification 336
- 9.4 Errors-in-Variables Bias 338
- 9.5 Sample Selection Bias 340
- 9.6 Simultaneous Causality Bias 343
- 9.7 Threats to the Internal Validity of a Multiple Regression Study 344

### **PART THREE Further Topics in Regression Analysis**

---

- 10.1 Notation for Panel Data 362
- 10.2 The Fixed Effects Regression Model 369
- 10.3 The Fixed Effects Regression Assumptions 375
- 11.1 The Linear Probability Model 396
- 11.2 The Probit Model, Predicted Probabilities, and Estimated Effects 400
- 11.3 Logit Regression 402
- 12.1 The General Instrumental Variables Regression Model and Terminology 438
- 12.2 Two Stage Least Squares 440
- 12.3 The Two Conditions for Valid Instruments 441
- 12.4 The IV Regression Assumptions 442
- 12.5 A Rule of Thumb for Checking for Weak Instruments 446
- 12.6 The Overidentifying Restrictions Test (The  $J$ -Statistic) 449
- 14.1  $m$ -Fold Cross Validation 523
- 14.2 The Principal Components of  $X$  535

### **PART FOUR Regression Analysis of Economic Time Series Data**

---

- 15.1 Lags, First Differences, Logarithms, and Growth Rates 557
- 15.2 Autocorrelation (Serial Correlation) and Autocovariance 559
- 15.3 Stationarity 562
- 15.4 Autoregressions 568
- 15.5 The Autoregressive Distributed Lag Model 571
- 15.6 The Least Squares Assumptions for Forecasting with Time Series Data 572
- 15.7 Pseudo Out-of-Sample Forecasts 575
- 15.8 The QLR Test for Coefficient Stability 592
- 16.1 The Distributed Lag Model and Exogeneity 616

- 16.2 The Distributed Lag Model Assumptions 618
- 16.3 HAC Standard Errors 624
- 17.1 Vector Autoregressions 650
- 17.2 Iterated Multi-period Forecasts 656
- 17.3 Direct Multi-period Forecasts 658
- 17.4 Orders of Integration, Differencing, and Stationarity 660
- 17.5 Cointegration 664

## **PART FIVE**    **Regression Analysis of Economic Time Series Data**

---

- 18.1 The Extended Least Squares Assumptions for Regression with a Single Regressor 689
- 19.1 The Extended Least Squares Assumptions in the Multiple Regression Model 715
- 19.2 The Multivariate Central Limit Theorem 718
- 19.3 Gauss–Markov Theorem for Multiple Regression 727
- 19.4 The GLS Assumptions 729

This page intentionally left blank

# General Interest Boxes

The Distribution of Adulthood Earnings in the United Kingdom by Childhood Socioeconomic Circumstances	72
The Unpegging of the Swiss Franc	77
Financial Diversification and Portfolios	84
Off the Mark!	108
Social Class or Education? Childhood Circumstances and Adult Earnings Revisited	122
A Way to Increase Voter Turnout	124
The “Beta” of a Stock	152
The Economic Value of a Year of Education: Homoskedasticity or Heteroskedasticity?	193
Is Coffee Good for Your Health?	214
The Effect of Ageing on Healthcare Expenditures: A Red Herring?	304
The Demand for Economics Journals	307
Do Stock Mutual Funds Outperform the Market?	341
James Heckman and Daniel McFadden, Nobel Laureates	414
When Was Instrumental Variables Regression Invented?	430
The First IV Regression	447
The Externalities of Smoking	451
The Hawthorne Effect	480
Conditional Cash Transfers in Rural Mexico to Increase School Enrollment	483
Text as Data	543
Can You Beat the Market?	564
The River of Blood	577
Orange Trees on the March	635
NEWS FLASH: Commodity Traders Send Shivers Through Disney World	636
Nobel Laureates in Time Series Econometrics	680

This page intentionally left blank

# Preface

Econometrics can be a fun course for both teacher and student. The real world of economics, business, and government is a complicated and messy place, full of competing ideas and questions that demand answers. Does healthcare spending actually improve health outcomes? Can you make money in the stock market by buying when prices are historically low, relative to earnings, or should you just sit tight, as the random walk theory of stock prices suggests? Does heavy intake of coffee lower the risk of disease or death? Econometrics helps us sort out sound ideas from crazy ones and find quantitative answers to important quantitative questions. Econometrics opens a window on our complicated world that lets us see the relationships on which people, businesses, and governments base their decisions.

*Introduction to Econometrics* is designed for a first course in undergraduate econometrics. It is our experience that to make econometrics relevant in an introductory course, interesting applications must motivate the theory and the theory must match the applications. This simple principle represents a significant departure from the older generation of econometrics books, in which theoretical models and assumptions do not match the applications. It is no wonder that some students question the relevance of econometrics after they spend much of their time learning assumptions that they subsequently realize are unrealistic so that they must then learn “solutions” to “problems” that arise when the applications do not match the assumptions. We believe that it is far better to motivate the need for tools with a concrete application and then to provide a few simple assumptions that match the application. Because the methods are immediately relevant to the applications, this approach can make econometrics come alive.

To improve student results, we recommend pairing the text content with MyLab Economics, which is the teaching and learning platform that empowers you to reach every student. By combining trusted author content with digital tools and a flexible platform, MyLab personalizes the learning experience and will help your students learn and retain key course concepts while developing skills that future employers are seeking in their candidates. MyLab Economics helps you teach your course, your way. Learn more at [www.pearson.com/mylab/economics](http://www.pearson.com/mylab/economics).

## New To This Edition

- New chapter on “Big Data” and machine learning
- Forecasting in time series data with large data sets

- Dynamic factor models
- Parallel treatment of prediction and causal inference using regression
- Coverage of realized volatility as well as autoregressive conditional heteroskedasticity
- Updated discussion of weak instruments

Very large data sets are increasingly being used in economics and related fields. Applications include predicting consumer choices, measuring the quality of hospitals or schools, analyzing nonstandard data such as text data, and macroeconomic forecasting with many variables. The three main additions in this edition incorporate the fundamentals of this growing and exciting area of application.

First, we have a new chapter (Chapter 14) that focuses on big data and machine learning methods. Within economics, many of the applications to date have focused on the so called many-predictor problem, where the number of predictors is large relative to the sample size—perhaps even exceeding the sample size. With many predictors, ordinary least squares (OLS) provides poor predictions, and other methods, such as the LASSO, can have much lower out-of-sample prediction errors. This chapter goes over the concepts of out-of-sample prediction, why OLS performs poorly, and how shrinkage can improve upon OLS. The chapter introduces shrinkage methods and prediction using principal components, shows how to choose tuning parameters by cross-validation, and explains how these methods can be used to analyze nonstandard data such as text data. As usual, this chapter has a running empirical example, in this case, prediction of school-level test scores given school-level characteristics, for California elementary schools.

Second, in Chapter 17 (newly renumbered), we extend the many-predictor focus of Chapter 14 to time series data. Specifically, we show how the dynamic factor model can handle a very large number of time series, and show how to implement the dynamic factor model using principal components analysis. We illustrate the dynamic factor model and its use for forecasting with a 131-variable dataset of U.S. quarterly macroeconomic time series.

Third, we now lay out these two uses of regression—causal inference and prediction—up front, when regression is first introduced in Chapter 4. Regression is a statistical tool that can be used to make causal inferences or to make predictions; the two applications place different demands on how the data are collected. When the data are from a randomized controlled experiment, OLS estimates the causal effect. In observational data, if we are interested in estimating the causal effect, then the econometrician needs to use control variables and/or instruments to produce as-if randomization of the variable of interest. In contrast, for prediction, one is not interested in the causal effect so one does not need as-if random variation; however, the estimation (“training”) data set must be drawn from the same population as the observations for which one wishes to make the prediction.

This edition has several smaller changes. For example, we now introduce realized volatility as a complement to the GARCH model when analyzing time series data with volatility clustering. In addition, we now extend the discussion (in a new general interest box) of the historical origins of instrumental variables regression in Chapter 12. This treatment now includes a first-ever reproduction of the original derivation of the IV estimator, which was in a letter from Philip Wright to his son Sewall in the spring of 1926, and a discussion of the first IV regression, an estimate of the elasticity of supply of flaxseed.

## Solving Teaching and Learning Challenges

*Introduction to Econometrics* differs from other texts in three main ways. First, we integrate real-world questions and data into the development of the theory, and we take seriously the substantive findings of the resulting empirical analysis. Second, our choice of topics reflects modern theory and practice. Third, we provide theory and assumptions that match the applications. Our aim is to teach students to become sophisticated consumers of econometrics and to do so at a level of mathematics appropriate for an introductory course.

### Real-World Questions and Data

We organize each methodological topic around an important real-world question that demands a specific numerical answer. For example, we teach single-variable regression, multiple regression, and functional form analysis in the context of estimating the effect of school inputs on school outputs. (Do smaller elementary school class sizes produce higher test scores?) We teach panel data methods in the context of analyzing the effect of drunk driving laws on traffic fatalities. We use possible racial discrimination in the market for home loans as the empirical application for teaching regression with a binary dependent variable (logit and probit). We teach instrumental variable estimation in the context of estimating the demand elasticity for cigarettes. Although these examples involve economic reasoning, all can be understood with only a single introductory course in economics, and many can be understood without any previous economics coursework. Thus the instructor can focus on teaching econometrics, not microeconomics or macroeconomics.

We treat all our empirical applications seriously and in a way that shows students how they can learn from data but at the same time be self-critical and aware of the limitations of empirical analyses. Through each application, we teach students to explore alternative specifications and thereby to assess whether their substantive findings are robust. The questions asked in the empirical applications are important, and we provide serious and, we think, credible answers. We encourage students and instructors to disagree, however, and invite them to reanalyze the

data, which are provided on the text's Companion Website ([www.pearsonglobaleditions.com](http://www.pearsonglobaleditions.com)) and in MyLab Economics.

Throughout the text, we have focused on helping students understand, retain, and apply the essential ideas. *Chapter introductions* provide real-world grounding and motivation, as well as brief road maps highlighting the sequence of the discussion. *Key terms* are boldfaced and defined in context throughout each chapter, and *Key Concept boxes* at regular intervals recap the central ideas. *General interest boxes* provide interesting excursions into related topics and highlight real-world studies that use the methods or concepts being discussed in the text. A *Summary* concluding each chapter serves as a helpful framework for reviewing the main points of coverage.

Available for student practice or instructor assignment in MyLab Economics are *Review the Concepts questions*, *Exercises*, and *Empirical Exercises* from the text. These questions and exercises are auto-graded, giving students practical hands-on experience with solving problems using the data sets used in the text.

- 100 percent of Review the Concepts questions are available in MyLab.
- Select Exercises and Empirical Exercises are available in MyLab. Many of the Empirical Exercises are algorithmic and based on the data sets used in the text. These exercises require students to use Excel or an econometrics software package to analyze the data and derive results.
- New to the 4<sup>th</sup> edition are concept exercises that focus on core concepts and economic interpretations. Many are algorithmic and include the Help Me Solve This learning aid.

## Contemporary Choice of Topics

The topics we cover reflect the best of contemporary applied econometrics. One can only do so much in an introductory course, so we focus on procedures and tests that are commonly (or increasingly) used in practice. For example:

- ***Instrumental variables regression.*** We present instrumental variables regression as a general method for handling correlation between the error term and a regressor, which can arise for many reasons, including omitted variables and simultaneous causality. The two assumptions for a valid instrument—exogeneity and relevance—are given equal billing. We follow that presentation with an extended discussion of where instruments come from and with tests of overidentifying restrictions and diagnostics for weak instruments, and we explain what to do if these diagnostics suggest problems.
- ***Program evaluation.*** Many modern econometric studies analyze either randomized controlled experiments or quasi-experiments, also known as natural experiments. We address these topics, often collectively referred to as program

evaluation, in Chapter 13. We present this research strategy as an alternative approach to the problems of omitted variables, simultaneous causality, and selection, and we assess both the strengths and the weaknesses of studies using experimental or quasi-experimental data.

- **Prediction with “big data.”** Chapter 14 takes up the opportunities and challenges posed by large cross-sectional data sets. An increasingly common application in econometrics is making predictions when the number of predictors is very large. This chapter focuses on methods designed to use many predictors in a way that produces accurate and precise out-of-sample predictions. The chapter covers some of the building blocks of machine learning, and the methods can substantially improve upon OLS when the number of predictors is large. In addition, these methods extend to nonstandard data, such as text data.
- **Forecasting.** The chapter on forecasting (Chapter 15) considers univariate (autoregressive) and multivariate forecasts using time series regression, not large simultaneous equation structural models. We focus on simple and reliable tools, such as autoregressions and model selection via an information criterion, that work well in practice. This chapter also features a practically oriented treatment of structural breaks (at known and unknown dates) and pseudo out-of-sample forecasting, all in the context of developing stable and reliable time series forecasting models.
- **Time series regression.** The chapter on causal inference using time series data (Chapter 16) pays careful attention to when different estimation methods, including generalized least squares, will or will not lead to valid causal inferences and when it is advisable to estimate dynamic regressions using OLS with heteroskedasticity- and autocorrelation-consistent standard errors.

## Theory That Matches Applications

Although econometric tools are best motivated by empirical applications, students need to learn enough econometric theory to understand the strengths and limitations of those tools. We provide a modern treatment in which the fit between theory and applications is as tight as possible, while keeping the mathematics at a level that requires only algebra.

Modern empirical applications share some common characteristics: The data sets typically have many observations (hundreds or more); regressors are not fixed over repeated samples but rather are collected by random sampling (or some other mechanism that makes them random); the data are not normally distributed; and there is no *a priori* reason to think that the errors are homoskedastic (although often there are reasons to think that they are heteroskedastic).

These observations lead to important differences between the theoretical development in this text and other texts:

- **Large-sample approach.** Because data sets are large, from the outset we use large-sample normal approximations to sampling distributions for hypothesis testing and confidence intervals. In our experience, it takes less time to teach the rudiments of large-sample approximations than to teach the Student  $t$  and exact  $F$  distributions, degrees-of-freedom corrections, and so forth. This large-sample approach also saves students the frustration of discovering that, because of nonnormal errors, the exact distribution theory they just mastered is irrelevant. Once taught in the context of the sample mean, the large-sample approach to hypothesis testing and confidence intervals carries directly through multiple regression analysis, logit and probit, instrumental variables estimation, and time series methods.
- **Random sampling.** Because regressors are rarely fixed in econometric applications, from the outset we treat data on all variables (dependent and independent) as the result of random sampling. This assumption matches our initial applications to cross-sectional data, it extends readily to panel and time series data, and because of our large-sample approach, it poses no additional conceptual or mathematical difficulties.
- **Heteroskedasticity.** Applied econometricians routinely use heteroskedasticity-robust standard errors to eliminate worries about whether heteroskedasticity is present or not. In this book, we move beyond treating heteroskedasticity as an exception or a “problem” to be “solved”; instead, we allow for heteroskedasticity from the outset and simply use heteroskedasticity-robust standard errors. We present homoskedasticity as a special case that provides a theoretical motivation for OLS.

## Skilled Producers, Sophisticated Consumers

We hope that students using this book will become sophisticated consumers of empirical analysis. To do so, they must learn not only how to use the tools of regression analysis but also how to assess the validity of empirical analyses presented to them.

Our approach to teaching how to assess an empirical study is threefold. First, immediately after introducing the main tools of regression analysis, we devote Chapter 9 to the threats to internal and external validity of an empirical study. This chapter discusses data problems and issues of generalizing findings to other settings. It also examines the main threats to regression analysis, including omitted variables, functional form misspecification, errors-in-variables, selection, and simultaneity—and ways to recognize these threats in practice.

Second, we apply these methods for assessing empirical studies to the empirical analysis of the ongoing examples in the book. We do so by considering alternative specifications and by systematically addressing the various threats to validity of the analyses presented in the book.

Third, to become sophisticated consumers, students need firsthand experience as producers. Active learning beats passive learning, and econometrics is an ideal course for active learning. For this reason, the MyLab Economics and text website feature data sets, software, and suggestions for empirical exercises of different scopes.

## Approach to Mathematics and Level of Rigor

Our aim is for students to develop a sophisticated understanding of the tools of modern regression analysis, whether the course is taught at a “high” or a “low” level of mathematics. Parts I through IV of the text (which cover the substantive material) are written for students with only precalculus mathematics. Parts I through IV have fewer equations and more applications than many introductory econometrics books and far fewer equations than books aimed at mathematical sections of undergraduate courses. But more equations do not imply a more sophisticated treatment. In our experience, a more mathematical treatment does not lead to a deeper understanding for most students.

That said, different students learn differently, and for mathematically well-prepared students, learning can be enhanced by a more explicit mathematical treatment. The appendices in Parts I-IV therefore provide key calculations that are too involved to be included in the text. In addition, Part V contains an introduction to econometric theory that is appropriate for students with a stronger mathematical background. When the mathematical chapters in Part V are used in conjunction with the material in Parts I through IV (including appendices), this book is suitable for advanced undergraduate or master’s level econometrics courses.

## Developing Career Skills

For students to succeed in a rapidly changing job market, they should be aware of their career options and how to go about developing a variety of skills. Data analysis is an increasingly marketable skill. This text prepares the students for a range of data analytic applications, including causal inference and prediction. It also introduces the students to the core concepts of prediction using large data sets.

## Table of Contents Overview

There are five parts to *Introduction to Econometrics*. This text assumes that the student has had a course in probability and statistics, although we review that material in Part I. We cover the core material of regression analysis in Part II. Parts III, IV, and V present additional topics that build on the core treatment in Part II.

### Part I

Chapter 1 introduces econometrics and stresses the importance of providing quantitative answers to quantitative questions. It discusses the concept of causality in statistical studies and surveys the different types of data encountered in econometrics. Material from probability and statistics is reviewed in Chapters 2 and 3, respectively; whether these chapters are taught in a given course or are simply provided as a reference depends on the background of the students.

### Part II

Chapter 4 introduces regression with a single regressor and ordinary least squares (OLS) estimation, and Chapter 5 discusses hypothesis tests and confidence intervals in the regression model with a single regressor. In Chapter 6, students learn how they can address omitted variable bias using multiple regression, thereby estimating the effect of one independent variable while holding other independent variables constant. Chapter 7 covers hypothesis tests, including  $F$ -tests, and confidence intervals in multiple regression. In Chapter 8, the linear regression model is extended to models with nonlinear population regression functions, with a focus on regression functions that are linear in the parameters (so that the parameters can be estimated by OLS). In Chapter 9, students step back and learn how to identify the strengths and limitations of regression studies, seeing in the process how to apply the concepts of internal and external validity.

### Part III

Part III presents extensions of regression methods. In Chapter 10, students learn how to use panel data to control for unobserved variables that are constant over time. Chapter 11 covers regression with a binary dependent variable. Chapter 12 shows how instrumental variables regression can be used to address a variety of problems that produce correlation between the error term and the regressor, and examines how one might find and evaluate valid instruments. Chapter 13 introduces students to the analysis of data from experiments and quasi-, or natural, experiments, topics often referred to as “program evaluation.” Chapter 14 turns to econometric issues that arise with large data sets, and focuses on prediction when there are very many predictors.

## Part IV

Part IV takes up regression with time series data. Chapter 15 focuses on forecasting and introduces various modern tools for analyzing time series regressions, such as tests for stability. Chapter 16 discusses the use of time series data to estimate causal relations. Chapter 17 presents some more advanced tools for time series analysis, including models of volatility clustering and dynamic factor models.

## Part V

Part V is an introduction to econometric theory. This part is more than an appendix that fills in mathematical details omitted from the text. Rather, it is a self-contained treatment of the econometric theory of estimation and inference in the linear regression model. Chapter 18 develops the theory of regression analysis for a single regressor; the exposition does not use matrix algebra, although it does demand a higher level of mathematical sophistication than the rest of the text. Chapter 19 presents the multiple regression model, instrumental variables regression, generalized method of moments estimation of the linear model, and principal components analysis, all in matrix form.

## Prerequisites Within the Book

Because different instructors like to emphasize different material, we wrote this book with diverse teaching preferences in mind. To the maximum extent possible, the chapters in Parts III, IV, and V are “stand-alone” in the sense that they do not require first teaching all the preceding chapters. The specific prerequisites for each chapter are described in Table I. Although we have found that the sequence of topics adopted in the text works well in our own courses, the chapters are written in a way that allows instructors to present topics in a different order if they so desire.

## Sample Courses

This book accommodates several different course structures.

**TABLE I** Guide to Prerequisites for Special-Topic Chapters in Parts III, IV, and V

Chapter	Prerequisite parts or chapters								
	Part I	Part II		Part III		Part IV			Part V
	1-3	4-7, 9	8	10.1, 10.2	12.1, 12.2	15.1-15.4	15.5-15.8	16	18
10	X <sup>a</sup>	X <sup>a</sup>	X						
11	X <sup>a</sup>	X <sup>a</sup>	X						
12.1, 12.2	X <sup>a</sup>	X <sup>a</sup>	X						
12.3-12.6	X <sup>a</sup>	X <sup>a</sup>	X	X	X				
13	X <sup>a</sup>	X <sup>a</sup>	X	X	X				
14	X <sup>a</sup>	X <sup>a</sup>	X						
15	X <sup>a</sup>	X <sup>a</sup>	b						
16	X <sup>a</sup>	X <sup>a</sup>	b			X			
17	X <sup>a</sup>	X <sup>a</sup>	b			X	X	X	
18	X	X	X						
19	X	X	X		X				X

This table shows the minimum prerequisites needed to cover the material in a given chapter. For example, estimation of dynamic causal effects with time series data (Chapter 16) first requires Part I (as needed, depending on student preparation, and except as noted in footnote a), Part II (except for Chapter 8; see footnote b), and Sections 15.1 through 15.4.

<sup>a</sup>Chapters 10 through 17 use exclusively large-sample approximations to sampling distributions, so the optional Sections 3.6 (the Student  $t$  distribution for testing means) and 5.6 (the Student  $t$  distribution for testing regression coefficients) can be skipped.

<sup>b</sup>Chapters 15 through 17 (the time series chapters) can be taught without first teaching Chapter 8 (nonlinear regression functions) if the instructor pauses to explain the use of logarithmic transformations to approximate percentage changes.

### Standard Introductory Econometrics

This course introduces econometrics (Chapter 1) and reviews probability and statistics as needed (Chapters 2 and 3). It then moves on to regression with a single regressor, multiple regression, the basics of functional form analysis, and the evaluation of regression studies (all Part II). The course proceeds to cover regression with panel data (Chapter 10), regression with a limited dependent variable (Chapter 11), and instrumental variables regression (Chapter 12), as time permits. The course then

turns to experiments and quasi-experiments in Chapter 13, topics that provide an opportunity to return to the questions of estimating causal effects raised at the beginning of the semester and to recapitulate core regression methods. If there is time, the students can be introduced to big data and machine learning methods at the end (Chapter 14). *Prerequisites: Algebra II and introductory statistics.*

## Introductory Econometrics with Time Series and Forecasting Applications

Like a standard introductory course, this course covers all of Part I (as needed) and Part II. Optionally, the course next provides a brief introduction to panel data (Sections 10.1 and 10.2) and takes up instrumental variables regression (Chapter 12, or just Sections 12.1 and 12.2). The course then proceeds to Chapter 14 (prediction in large cross sectional data sets). It then turns to Part IV, covering forecasting (Chapter 15) and estimation of dynamic causal effects (Chapter 16). If time permits, the course can include some advanced topics in time series analysis such as volatility clustering (Section 17.5) and forecasting with many predictors (Section 17.6). *Prerequisites: Algebra II and introductory statistics.*

## Applied Time Series Analysis and Forecasting

This book also can be used for a short course on applied time series and forecasting, for which a course on regression analysis is a prerequisite. Some time is spent reviewing the tools of basic regression analysis in Part II, depending on student preparation. The course then moves directly to time series forecasting (Chapter 15), estimation of dynamic causal effects (Chapter 16), and advanced topics in time series analysis (Chapter 17), including vector autoregressions. If there is time, the course can cover prediction using large data sets (Chapter 14 and Section 17.6). An important component of this course is hands-on forecasting exercises, available as the end-of-chapter Empirical Exercises for Chapters 15 and 17. *Prerequisites: Algebra II and basic introductory econometrics or the equivalent.*

## Introduction to Econometric Theory

This book is also suitable for an advanced undergraduate course in which the students have a strong mathematical preparation or for a master's level course in econometrics. The course briefly reviews the theory of statistics and probability as necessary (Part I). The course introduces regression analysis using the nonmathematical, applications-based treatment of Part II. This introduction is followed by the theoretical development in Chapters 18 and 19 (through Section 19.5). The course then takes up regression with a limited dependent variable (Chapter 11) and maximum likelihood estimation (Appendix 11.2). Next, the course optionally turns to instrumental variables regression and generalized method of moments (Chapter 12 and Section 19.7), time series methods (Chapter 15), the estimation of

causal effects using time series data and generalized least squares (Chapter 16 and Section 19.6), and/or to machine learning methods (Chapter 14 and Appendix 19.7). *Prerequisites: Calculus and introductory statistics. Chapter 18 assumes previous exposure to matrix algebra.*

## Instructor Teaching Resources

This program comes with the following teaching resources:

Supplements available to instructors at <a href="http://www.pearsonglobaleditions.com">www.pearsonglobaleditions.com</a>	Features of the Supplement
<b>Solutions Manual</b>	Solutions to the end-of-chapter content.
<b>Test Bank</b> Authored by Manfred Keil, Claremont McKenna College	1,000 multiple-choice questions, essays and longer questions, and mathematical and graphical problems with these annotations: <ul style="list-style-type: none"> <li>Type (Multiple-choice, essay, graphical)</li> </ul>
<b>Computerized TestGen</b>	TestGen allows instructors to: <ul style="list-style-type: none"> <li>Customize, save, and generate classroom tests</li> <li>Edit, add, or delete questions from the Test Item Files</li> <li>Analyze test results</li> <li>Organize a database of tests and student results.</li> </ul>
<b>PowerPoints</b>	Slides include all the graphs, tables, and equations in the text.  PowerPoints meet accessibility standards for students with disabilities. Features include, but not limited to: <ul style="list-style-type: none"> <li>Keyboard and Screen Reader access</li> <li>Alternative text for images</li> <li>High color contrast between background and foreground colors</li> </ul>
<b>Companion Website</b>	The Companion Website provides a wide range of additional resources for students and faculty. These resources include more and more in depth empirical exercises, data sets for the empirical exercises, replication files for empirical results reported in the text, and EViews tutorials.

## Acknowledgments

A great many people contributed to the first edition of this book. Our biggest debts of gratitude are to our colleagues at Harvard and Princeton who used early drafts of this book in their classrooms. At Harvard's Kennedy School of Government, Suzanne Cooper provided invaluable suggestions and detailed comments on multiple drafts. As a coteacher with one of the authors (Stock), she also helped vet much of the material in this book while it was being developed for a required course for master's students at the Kennedy School. We are also indebted to two other Kennedy School colleagues at the time, Alberto Abadie and Sue Dynarski, for their patient explanations of quasi-experiments and the field of program evaluation and for their detailed comments on early drafts of the text. At Princeton, Eli Tamer taught from an early draft and also provided helpful comments on the penultimate draft of the book.

We also owe much to many of our friends and colleagues in econometrics who spent time talking with us about the substance of this book and who collectively made so many helpful suggestions. Bruce Hansen (University of Wisconsin–Madison) and Bo Honore (Princeton) provided helpful feedback on very early outlines and preliminary versions of the core material in Part II. Joshua Angrist (MIT) and Guido Imbens (University of California, Berkeley) provided thoughtful suggestions about our treatment of materials on program evaluation. Our presentation of the material on time series has benefited from discussions with Yacine Ait-Sahalia (Princeton), Graham Elliott (University of California, San Diego), Andrew Harvey (Cambridge University), and Christopher Sims (Princeton). Finally, many people made helpful suggestions on parts of the manuscript close to their area of expertise: Don Andrews (Yale), John Bound (University of Michigan), Gregory Chow (Princeton), Thomas Downes (Tufts), David Drukker (StataCorp.), Jean Baldwin Grossman (Princeton), Eric Hanushek (Hoover Institution), James Heckman (University of Chicago), Han Hong (Princeton), Caroline Hoxby (Harvard), Alan Krueger (Princeton), Steven Levitt (University of Chicago), Richard Light (Harvard), David Neumark (Michigan State University), Joseph Newhouse (Harvard), Pierre Perron (Boston University), Kenneth Warner (University of Michigan), and Richard Zeckhauser (Harvard).

Many people were very generous in providing us with data. The California test score data were constructed with the assistance of Les Axelrod of the Standards and Assessments Division, California Department of Education. We are grateful to Charlie DePascale, Student Assessment Services, Massachusetts Department of Education, for his help with aspects of the Massachusetts test score data set. Christopher Ruhm (University of North Carolina, Greensboro) graciously provided us with his data set on drunk driving laws and traffic fatalities. The research department at the Federal Reserve Bank of Boston deserves thanks for putting together its data on racial discrimination in mortgage lending; we particularly thank Geoffrey Tootell for providing us with the updated version of the data set we use in Chapter 9 and Lynn Browne for explaining its policy context. We thank Jonathan Gruber (MIT) for sharing his data on cigarette sales, which we analyze in Chapter 12, and

Alan Krueger (Princeton) for his help with the Tennessee STAR data that we analyze in Chapter 13.

We thank several people for carefully checking the page proof for errors. Kerry Griffin and Yair Listokin read the entire manuscript, and Andrew Fraker, Ori Heffetz, Amber Henry, Hong Li, Alessandro Tarozzi, and Matt Watson worked through several chapters.

In the first edition, we benefited from the help of an exceptional development editor, Jane Tufts, whose creativity, hard work, and attention to detail improved the book in many ways, large and small. Pearson provided us with first-rate support, starting with our excellent editor, Sylvia Mallory, and extending through the entire publishing team. Jane and Sylvia patiently taught us a lot about writing, organization, and presentation, and their efforts are evident on every page of this book. We extend our thanks to the superb Pearson team, who worked with us on the second edition: Adrienne D'Ambrosio (senior acquisitions editor), Bridget Page (associate media producer), Charles Spaulding (senior designer), Nancy Fenton (managing editor) and her selection of Nancy Freihofer and Thompson Steele Inc. who handled the entire production process, Heather McNally (supplements coordinator), and Denise Clinton (editor-in-chief). Finally, we had the benefit of Kay Ueno's skilled editing in the second edition. We are also grateful to the excellent third edition Pearson team of Adrienne D'Ambrosio, Nancy Fenton, and Jill Kolongowski, as well as Rose Kernan, the project manager with Cenveo Publisher Services. We also wish to thank the Pearson team who worked on the fourth edition: David Alexander, Christine Donovan, Nicole Nedwitek, and Rose Kernan, project manager with Cenveo Publisher Services.

We also received a great deal of help and suggestions from faculty, students, and researchers as we prepared the third edition and its update. The changes made in the third edition incorporate or reflect suggestions, corrections, comments, data, and help provided by a number of researchers and instructors: Donald Andrews (Yale University), Jushan Bai (Columbia), James Cobbe (Florida State University), Susan Dynarski (University of Michigan), Nicole Eichelberger (Texas Tech University), Boyd Fjeldsted (University of Utah), Martina Grunow, Daniel Hamermesh (University of Texas–Austin), Keisuke Hirano (University of Arizona), Bo Honore (Princeton University), Guido Imbens (Harvard University), Manfred Keil (Claremont McKenna College), David Laibson (Harvard University), David Lee (Princeton University), Brigitte Madrian (Harvard University), Jorge Marquez (University of Maryland), Karen Bennett Mathis (Florida Department of Citrus), Alan Mehlenbacher (University of Victoria), Ulrich Müller (Princeton University), Serena Ng (Columbia University), Harry Patrinos (World Bank), Zhuan Pei (Brandeis University), Peter Summers (Texas Tech University), Andrey Vasnov (University of Sydney), and Douglas Young (Montana State University). We also benefited from student input from F. Hoces dela Guardia and Carrie Wilson.

Thoughtful reviews for the third edition were prepared for Pearson by Steve DeLoach (Elon University), Jeffrey DeSimone (University of Texas at Arlington), Gary V. Engelhardt (Syracuse University), Luca Flabbi (Georgetown University), Steffen Habermalz (Northwestern University), Carolyn J. Heinrich (University of Wisconsin–Madison), Emma M. Iglesias-Vazquez (Michigan State University), Carlos Lamarche

(University of Oklahoma), Vicki A. McCracken (Washington State University), Claudiney M. Pereira (Tulane University), and John T. Warner (Clemson University). We also received very helpful input on draft revisions of Chapters 7 and 10 from John Berdell (DePaul University), Janet Kohlhase (University of Houston), Aprajit Mahajan (Stanford University), Xia Meng (Brandeis University), and Chan Shen (Georgetown University). We thank Christopher Stock for helping with the third edition cover.

In the fourth edition, we benefited from detailed comments on our prior treatment of causal analysis by Judea Pearl (UCLA) and Bryant Chen. Gary Chamberlain (Harvard), Guido Imbens (Stanford) and Jan Speiss (Stanford) provided thoughtful comments and guidance on Chapter 14. We received additional comments and/or corrections from Carlos C. Bautista (University of the Philippines), Brian Bethune (Tufts), Amitabh Chandra (Harvard Kennedy School), Julia Chang (University of New South Wales), Maia Güell (University of Edinburgh), Greg Mankiw (Harvard), Alan Mehlenbacher (University of Victoria), Franco Peracchi (Tor Vergata University), Peter Siminski (University of Wollongong), Jinhua Wang (University of Cambridge), and Michael Wolf (University of Zurich). We also benefited from a review panel that focused on the new Chapter 14, comprised of Chrystie Burr (University of Colorado-Boulder), Bentley Coffey (University of South Carolina), and Galin Todorov (Florida Atlantic University).

Above all, we are indebted to our families for their endurance throughout this project. Writing this book took a long time, and for them, the project must have seemed endless. They, more than anyone else, bore the burden of this commitment, and for their help and support we are deeply grateful.

## Global Acknowledgments

We would like to thank the people who have contributed towards developing this book for the global markets and who have put in effort to update this global edition for students across the world.

Samprit Chakrabarti, International School of Business and Media  
 Raghvi Garg, Ashoka University  
 Daniel Howdon, University of Leeds  
 James Lomas, The University of York  
 Anisha Sharma, University of Oxford

We would also like to thank the individuals who reviewed the text and whose feedback has made this a better book.

Mostafa AboElsoud, Suez Canal University  
 Martin Christopher Arnold, University of Duisburg-Essen  
 Chitrita Bhowmick Chakrabarti, Victoria Institution  
 Jose Olmo, University of Southampton  
 Dragos Radu, King's College London  
 Raymond Wong, The University of Hong Kong

This page intentionally left blank

Ask a half dozen econometricians what econometrics is, and you could get a half dozen different answers. One might tell you that econometrics is the science of testing economic theories. A second might tell you that econometrics is the set of tools used for forecasting future values of economic variables, such as a firm's sales, the overall growth of the economy, or stock prices. Another might say that econometrics is the process of fitting mathematical economic models to real-world data. A fourth might tell you that it is the science and art of using historical data to make numerical, or quantitative, policy recommendations in government and business.

In fact, all these answers are right. At a broad level, econometrics is the science and art of using economic theory and statistical techniques to analyze economic data. Econometric methods are used in many branches of economics, including finance, labor economics, macroeconomics, microeconomics, marketing, and economic policy. Econometric methods are also commonly used in other social sciences, including political science and sociology.

This text introduces you to the core set of methods used by econometricians. We will use these methods to answer a variety of specific, quantitative questions from the worlds of business and government policy. This chapter poses four of those questions and discusses, in general terms, the econometric approach to answering them. The chapter concludes with a survey of the main types of data available to econometricians for answering these and other quantitative economic questions.

## 1.1 Economic Questions We Examine

Many decisions in economics, business, and government hinge on understanding relationships among variables in the world around us. These decisions require quantitative answers to quantitative questions.

This text examines several quantitative questions taken from current issues in economics. Four of these questions concern education policy, racial bias in mortgage lending, cigarette consumption, and macroeconomic forecasting.

### Question #1: Does Reducing Class Size Improve Elementary School Education?

Proposals for reform of the U.S. public education system generate heated debate. Many of the proposals concern the youngest students, those in elementary schools. Elementary school education has various objectives, such as developing social skills,

but for many parents and educators, the most important objective is basic academic learning: reading, writing, and basic mathematics. One prominent proposal for improving basic learning is to reduce class sizes at elementary schools. With fewer students in the classroom, the argument goes, each student gets more of the teacher's attention, there are fewer class disruptions, learning is enhanced, and grades improve.

But what, precisely, is the effect on elementary school education of reducing class size? Reducing class size costs money: It requires hiring more teachers and, if the school is already at capacity, building more classrooms. A decision maker contemplating hiring more teachers must weigh these costs against the benefits. To weigh costs and benefits, however, the decision maker must have a precise quantitative understanding of the likely benefits. Is the beneficial effect on basic learning of smaller classes large or small? Is it possible that smaller class size actually has no effect on basic learning?

Although common sense and everyday experience may suggest that more learning occurs when there are fewer students, common sense cannot provide a quantitative answer to the question of what exactly is the effect on basic learning of reducing class size. To provide such an answer, we must examine empirical evidence—that is, evidence based on data—relating class size to basic learning in elementary schools.

In this text, we examine the relationship between class size and basic learning, using data gathered from 420 California school districts in 1999. In the California data, students in districts with small class sizes tend to perform better on standardized tests than students in districts with larger classes. While this fact is consistent with the idea that smaller classes produce better test scores, it might simply reflect many other advantages that students in districts with small classes have over their counterparts in districts with large classes. For example, districts with small class sizes tend to have wealthier residents than districts with large classes, so students in small-class districts could have more opportunities for learning outside the classroom. It could be these extra learning opportunities that lead to higher test scores, not smaller class sizes. In Part II, we use multiple regression analysis to isolate the effect of changes in class size from changes in other factors, such as the economic background of the students.

## Question #2: Is There Racial Discrimination in the Market for Home Loans?

Most people buy their homes with the help of a mortgage, a large loan secured by the value of the home. By law, U.S. lending institutions cannot take race into account when deciding to grant or deny a request for a mortgage: Applicants who are identical in all ways except their race should be equally likely to have their mortgage applications approved. In theory, then, there should be no racial bias in mortgage lending.

In contrast to this theoretical conclusion, researchers at the Federal Reserve Bank of Boston found (using data from the early 1990s) that 28% of black applicants are

denied mortgages, while only 9% of white applicants are denied. Do these data indicate that, in practice, there is racial bias in mortgage lending? If so, how large is it?

The fact that more black than white applicants are denied in the Boston Fed data does not by itself provide evidence of discrimination by mortgage lenders because the black and white applicants differ in many ways other than their race. Before concluding that there is bias in the mortgage market, these data must be examined more closely to see if there is a difference in the probability of being denied for *otherwise identical* applicants and, if so, whether this difference is large or small. To do so, in Chapter 11 we introduce econometric methods that make it possible to quantify the effect of race on the chance of obtaining a mortgage, *holding constant* other applicant characteristics, notably their ability to repay the loan.

### Question #3: Does Healthcare Spending Improve Health Outcomes?

It is self-evident that no one lives forever, but avoidable deaths can be reduced and survival can be extended through the provision of healthcare. Healthcare has other beneficial effects too, like the improvement of the health-related quality of life of individuals. To these ends and more, a vast quantity of resources is devoted to the provision of healthcare worldwide. What is more there is enormous variation in the healthcare expenditures across countries both in absolute and per capita terms, as well as variations in health outcomes across countries, for example measured by life expectancy at birth.

Putting aside concerns about *iatrogenesis* (the idea that healthcare is bad for your health), basic economics says that more expenditure on healthcare should generally reduce avoidable mortality. But by how much? If the amount spent on healthcare increases by 1%, by what percentage will avoidable mortality decrease? The percentage change in avoidable mortality resulting from a 1% increase in healthcare expenditure is the *spending elasticity for mortality* (analogous to the *price elasticity of demand*, which is the percentage change in quantity demanded from a 1% increase in price). If we want to reduce avoidable mortality, say, 20% by increasing healthcare expenditure, then we need to know the spending elasticity for mortality to calculate the healthcare expenditure increase necessary to achieve this reduction in avoidable mortality.

A number of policy objectives are based on meeting targets based on avoidable mortality; for example, one of the United Nations Development Programme's sustainable development goals is that all countries should aim to reduce "under-5 mortality to at least as low as 25 per 1,000 live births."<sup>1</sup> But how should the goal be met: from expanding healthcare services or other services? And if increasing healthcare spending is to form part of the mix of policies, by how much will it need to increase? The answers to these can be obtained with estimates of the spending elasticity for mortality.

<sup>1</sup>United Nations Development Programme (UNDP), *The Sustainable Development Goals (SDGs): Goal 3: Good health and well-being*, 2017.

While economic theory, such as the production function for health, helps us analyze the mix of inputs that may lead to improved health outcomes, it does not tell us the actual values for parameters such as the spending elasticity for mortality. To estimate the value, we must examine empirical evidence about the returns to health-care spending—either based on variations in spending across countries or within countries over time (or both). In other words, we need to analyze the data on how health outcomes and healthcare expenditures are related.

For many years economists have attempted to address this question by considering the data on healthcare expenditures and mortality rates across countries, but such empirical research is fraught with challenges. Two of the biggest challenges concern the extensive heterogeneity across countries. The first challenge is observable heterogeneity, which concerns factors that affect countries' mortality rates that may also be associated with healthcare expenditure, for example, the income per capita of each country. This can be controlled for using multiple regression analysis, as described in Part II, since these factors are observable to the analyst. The second and more troublesome challenge is the presence of unobservable heterogeneity. Unobserved factors may be important in the underlying processes determining both how decisions are made on how much money is spent on healthcare, and how the overall level of health outcome that is attained. These factors result in causality running in both directions—healthcare reduces mortality, but higher healthcare expenditure might be a response to unobserved factors, such as small natural disasters that increase mortality. Methods for handling this “simultaneous causality” are described in Chapter 12, applied to the different but conceptually similar context of estimating the price elasticity of cigarette demand.

#### Question #4: By How Much Will U.S. GDP Grow Next Year?

It seems that people always want a sneak preview of the future. What will sales be next year at a firm that is considering investing in new equipment? Will the stock market go up next month, and, if it does, by how much? Will city tax receipts next year cover planned expenditures on city services? Will your microeconomics exam next week focus on externalities or monopolies? Will Saturday be a nice day to go to the beach?

One aspect of the future in which macroeconomists are particularly interested is the growth of real economic activity, as measured by real gross domestic product (GDP), during the next year. A management consulting firm might advise a manufacturing client to expand its capacity based on an upbeat forecast of economic growth. Economists at the Federal Reserve Board in Washington, D.C., are mandated to set policy to keep real GDP near its potential in order to maximize employment. If they forecast anemic GDP growth over the next year, they might expand liquidity in the economy by reducing interest rates or other measures, in an attempt to boost economic activity.

Professional economists who rely on numerical forecasts use econometric models to make those forecasts. A forecaster's job is to predict the future by using the

past, and econometricians do this by using economic theory and statistical techniques to quantify relationships in historical data.

The data we use to forecast the growth rate of GDP include past values of GDP and the so-called term spread in the United States. The *term spread* is the difference between long-term and short-term interest rates. It measures, among other things, whether investors expect short-term interest rates to rise or fall in the future. The term spread is usually positive, but it tends to fall sharply before the onset of a recession. One of the GDP growth rate forecasts we develop and evaluate in Chapter 15 is based on the term spread.

## Quantitative Questions, Quantitative Answers

Each of these four questions requires a numerical answer. Economic theory provides clues about that answer—for example, cigarette consumption ought to go down when the price goes up—but the actual value of the number must be learned empirically, that is, by analyzing data. Because we use data to answer quantitative questions, our answers always have some uncertainty: A different set of data would produce a different numerical answer. Therefore, the conceptual framework for the analysis needs to provide both a numerical answer to the question and a measure of how precise the answer is.

The conceptual framework used in this text is the multiple regression model, the mainstay of econometrics. This model, introduced in Part II, provides a mathematical way to quantify how a change in one variable affects another variable, holding other things constant. For example, what effect does a change in class size have on test scores, *holding constant* or *controlling for* student characteristics (such as family income) that a school district administrator cannot control? What effect does your race have on your chances of having a mortgage application granted, *holding constant* other factors such as your ability to repay the loan? What effect does a 1% increase in the price of cigarettes have on cigarette consumption, *holding constant* the income of smokers and potential smokers? The multiple regression model and its extensions provide a framework for answering these questions using data and for quantifying the uncertainty associated with those answers.

## 1.2 Causal Effects and Idealized Experiments

Like many other questions encountered in econometrics, the first three questions in Section 1.1 concern causal relationships among variables. In common usage, an action is said to cause an outcome if the outcome is the direct result, or consequence, of that action. Touching a hot stove causes you to get burned, drinking water causes you to be less thirsty, putting air in your tires causes them to inflate, putting fertilizer on your tomato plants causes them to produce more tomatoes. Causality means that a specific action (applying fertilizer) leads to a specific, measurable consequence (more tomatoes).

## Estimation of Causal Effects

How best might we measure the causal effect on tomato yield (measured in kilograms) of applying a certain amount of fertilizer, say, 100 grams of fertilizer per square meter?

One way to measure this causal effect is to conduct an experiment. In that experiment, a horticultural researcher plants many plots of tomatoes. Each plot is tended identically, with one exception: Some plots get 100 grams of fertilizer per square meter, while the rest get none. Whether or not a plot is fertilized is determined randomly by a computer, ensuring that any other differences between the plots are unrelated to whether they receive fertilizer. At the end of the growing season, the horticulturalist weighs the harvest from each plot. The difference between the average yield per square meter of the treated and untreated plots is the effect on tomato production of the fertilizer treatment.

This is an example of a **randomized controlled experiment**. It is controlled in the sense that there are both a **control group** that receives no treatment (no fertilizer) and a **treatment group** that receives the treatment (100 g/m<sup>2</sup> of fertilizer). It is randomized in the sense that the treatment is assigned randomly. This random assignment eliminates the possibility of a systematic relationship between, for example, how sunny the plot is and whether it receives fertilizer so that the only systematic difference between the treatment and control groups is the treatment. If this experiment is properly implemented on a large enough scale, then it will yield an estimate of the causal effect on the outcome of interest (tomato production) of the treatment (applying 100 g/m<sup>2</sup> of fertilizer).

In this text, the **causal effect** is defined to be the effect on an outcome of a given action or treatment, as measured in an ideal randomized controlled experiment. In such an experiment, the only systematic reason for differences in outcomes between the treatment and control groups is the treatment itself.

It is possible to imagine an ideal randomized controlled experiment to answer each of the first three questions in Section 1.1. For example, to study class size, one can imagine randomly assigning “treatments” of different class sizes to different groups of students. If the experiment is designed and executed so that the only systematic difference between the groups of students is their class size, then in theory this experiment would estimate the effect on test scores of reducing class size, holding all else constant.

Experiments are used increasingly widely in econometrics. In many applications, however, they are not an option because they are unethical, impossible to execute satisfactorily, too time-consuming, or prohibitively expensive. Even with non-experimental data, the concept of an ideal randomized controlled experiment is important because it provides a definition of a causal effect.

## Prediction, Forecasting, and Causality

Although the first three questions in Section 1.1, concern causal effects, the fourth—forecasting the growth rate of GDP—does not.

Forecasting is a special case of what statisticians and econometricians call **prediction**, which is using information on some variables to make a statement about the value of another variable. A **forecast** is a prediction about the value of a variable in the future, like GDP growth next year.

You do not need to know a causal relationship to make a good prediction. A good way to “predict” whether it is raining is to observe whether pedestrians are using umbrellas, but the act of using an umbrella does not cause it to rain.

When one has a small number of predictors and the data do not evolve over time, the multiple regression methods of Part II can provide reliable predictions. Predictions can often be improved, however, if there is a large number of candidate predictors. Methods for using many predictors are covered in Chapter 14.

Forecasts—that is, predictions about the future—use data on variables that evolve over time, which introduces new challenges and opportunities. As we will see in Chapter 15, multiple regression analysis allows us to quantify historical relationships, to check whether those relationships have been stable over time, to make quantitative forecasts about the future, and to assess the accuracy of those forecasts.

## 1.3 Data: Sources and Types

In econometrics, data come from one of two sources: experiments or nonexperimental observations of the world. This text examines both experimental and nonexperimental data sets.

### Experimental versus Observational Data

**Experimental data** come from experiments designed to evaluate a treatment or policy or to investigate a causal effect. For example, the state of Tennessee financed a large randomized controlled experiment examining class size in the 1980s. In that experiment, which we examine in Chapter 13, thousands of students were randomly assigned to classes of different sizes for several years and were given standardized tests annually.

The Tennessee class size experiment cost millions of dollars and required the ongoing cooperation of many administrators, parents, and teachers over several years. Because real-world experiments with human subjects are difficult to administer and to control, they have flaws relative to ideal randomized controlled experiments. Moreover, in some circumstances, experiments are not only expensive and difficult to administer but also unethical. (Would it be ethical to offer randomly selected teenagers inexpensive cigarettes to see how many they buy?) Because of these financial, practical, and ethical problems, experiments in economics are relatively rare. Instead, most economic data are obtained by observing real-world behavior.

Data obtained by observing actual behavior outside an experimental setting are called **observational data**. Observational data are collected using surveys, such as telephone surveys of consumers, and administrative records, such as historical records on mortgage applications maintained by lending institutions.

Observational data pose major challenges to econometric attempts to estimate causal effects, and the tools of econometrics are designed to tackle these challenges. In the real world, levels of “treatment” (the amount of fertilizer in the tomato example, the student–teacher ratio in the class size example) are not assigned at random, so it is difficult to sort out the effect of the “treatment” from other relevant factors. Much of econometrics, and much of this text, is devoted to methods for meeting the challenges encountered when real-world data are used to estimate causal effects.

Whether the data are experimental or observational, data sets come in three main types: cross-sectional data, time series data, and panel data. In this text, you will encounter all three types.

### Cross-Sectional Data

Data on different entities—workers, consumers, firms, governmental units, and so forth—for a single time period are called **cross-sectional data**. For example, the data on test scores in California school districts are cross sectional. Those data are for 420 entities (school districts) for a single time period (1999). In general, the number of entities on which we have observations is denoted  $n$ ; so, for example, in the California data set,  $n = 420$ .

The California test score data set contains measurements of several different variables for each district. Some of these data are tabulated in Table 1.1. Each row lists data for a different district. For example, the average test score for the first district (“district 1”) is 690.8; this is the average of the math and science test scores for all fifth-graders in that district in 1999 on a standardized test (the Stanford Achievement Test). The average student–teacher ratio in that district is 17.89; that is, the number of students in district 1 divided by the number of classroom teachers in district 1

**TABLE 1.1** Selected Observations on Test Scores and Other Variables for California School Districts in 1999

Observation (District) Number	District Average Test Score (fifth grade)	Student-Teacher Ratio	Expenditure per Pupil (\$)	Percentage of Students Learning English
1	690.8	17.89	\$6385	0.0%
2	661.2	21.52	5099	4.6
3	643.6	18.70	5502	30.0
4	647.7	17.36	7102	0.0
5	640.8	18.67	5236	13.9
⋮	⋮	⋮	⋮	⋮
418	645.0	21.89	4403	24.3
419	672.2	20.20	4776	3.0
420	655.8	19.04	5993	5.0

*Note:* The California test score data set is described in Appendix 4.1.

is 1789. Average expenditure per pupil in district 1 is \$6385. The percentage of students in that district still learning English—that is, the percentage of students for whom English is a second language and who are not yet proficient in English—is 0%.

The remaining rows present data for other districts. The order of the rows is arbitrary, and the number of the district, which is called the **observation number**, is an arbitrarily assigned number that organizes the data. As you can see in the table, all the variables listed vary considerably.

With cross-sectional data, we can learn about relationships among variables by studying differences across people, firms, or other economic entities during a single time period.

## Time Series Data

**Time series data** are data for a single entity (person, firm, country) collected at multiple time periods. Our data set on the growth rate of GDP and the term spread in the United States is an example of a time series data set. The data set contains observations on two variables (the growth rate of GDP and the term spread) for a single entity (the United States) for 232 time periods. Each time period in this data set is a quarter of a year (the first quarter is January, February, and March; the second quarter is April, May, and June; and so forth). The observations in this data set begin in the first quarter of 1960, which is denoted 1960:Q1, and end in the fourth quarter of 2017 (2017:Q4). The number of observations (that is, time periods) in a time series data set is denoted  $T$ . Because there are 232 quarters from 1960:Q1 to 2017:Q4, this data set contains  $T = 232$  observations.

Some observations in this data set are listed in Table 1.2. The data in each row correspond to a different time period (year and quarter). In the first quarter of 1960,

**TABLE 1.2** Selected Observations on the Growth Rate of GDP and the Term Spread in the United States: Quarterly Data, 1960:Q1–2017:Q4

Observation Number	Date (year: quarter)	GDP Growth Rate (% at an annual rate)	Term Spread (percentage points)
1	1960:Q1	8.8%	0.6
2	1960:Q2	−1.5	1.3
3	1960:Q3	1.0	1.5
4	1960:Q4	−4.9	1.6
5	1961:Q1	2.7	1.4
⋮	⋮	⋮	⋮
230	2017:Q2	3.0	1.4
231	2017:Q3	3.1	1.2
232	2017:Q4	2.5	1.2

*Note:* The United States GDP and term spread data set is described in Appendix 15.1.

for example, GDP grew 8.8% at an annual rate. In other words, if GDP had continued growing for four quarters at its rate during the first quarter of 1960, the level of GDP would have increased by 8.8%. In the first quarter of 1960, the long-term interest rate was 4.5%, and the short-term interest rate was 3.9%; so their difference, the term spread, was 0.6 percentage points.

By tracking a single entity over time, time series data can be used to study the evolution of variables over time and to forecast future values of those variables.

## Panel Data

**Panel data**, also called **longitudinal data**, are data for multiple entities in which each entity is observed at two or more time periods. Our data on cigarette consumption and prices are an example of a panel data set, and selected variables and observations in that data set are listed in Table 1.3. The number of entities in a panel data set is denoted  $n$ , and the number of time periods is denoted  $T$ . In the cigarette data set, we have observations on  $n = 48$  continental U.S. states (entities) for  $T = 11$  years (time periods) from 1985 to 1995. Thus, there is a total of  $n \times T = 48 \times 11 = 528$  observations.

Some data from the cigarette consumption data set are listed in Table 1.3. The first block of 48 observations lists the data for each state in 1985, organized alphabetically from Alabama to Wyoming. The next block of 48 observations lists the data for

**TABLE 1.3** Selected Observations on Cigarette Sales, Prices, and Taxes, by State and Year for U.S. States, 1985–1995

Observation Number	State	Year	Cigarette Sales (packs per capita)	Average Price per Pack (including taxes)	Total Taxes (cigarette excise tax + sales tax)
1	Alabama	1985	116.5	\$1.022	\$0.333
2	Arkansas	1985	128.5	1.015	0.370
3	Arizona	1985	104.5	1.086	0.362
⋮	⋮	⋮	⋮	⋮	⋮
47	West Virginia	1985	112.8	1.089	0.382
48	Wyoming	1985	129.4	0.935	0.240
49	Alabama	1986	117.2	1.080	0.334
⋮	⋮	⋮	⋮	⋮	⋮
96	Wyoming	1986	127.8	1.007	0.240
97	Alabama	1987	115.8	1.135	0.335
⋮	⋮	⋮	⋮	⋮	⋮
528	Wyoming	1995	112.2	1.585	0.360

*Note:* The cigarette consumption data set is described in Appendix 12.1.

## Cross-Sectional, Time Series, and Panel Data

**KEY CONCEPT****1.1**

- Cross-sectional data consist of multiple entities observed at a single time period.
- Time series data consist of a single entity observed at multiple time periods.
- Panel data (also known as longitudinal data) consist of multiple entities, where each entity is observed at two or more time periods.

1986, and so forth, through 1995. For example, in 1985, cigarette sales in Arkansas were 128.5 packs per capita (the total number of packs of cigarettes sold in Arkansas in 1985 divided by the total population of Arkansas in 1985 equals 128.5). The average price of a pack of cigarettes in Arkansas in 1985, including tax, was \$1.015, of which 37¢ went to federal, state, and local taxes.

Panel data can be used to learn about economic relationships from the experiences of the many different entities in the data set and from the evolution over time of the variables for each entity.

The definitions of cross-sectional data, time series data, and panel data are summarized in Key Concept 1.1.

### Summary

1. Many decisions in business and economics require quantitative estimates of how a change in one variable affects another variable.
2. Conceptually, the way to estimate a causal effect is in an ideal randomized controlled experiment, but performing experiments in economic applications can be unethical, impractical, or too expensive.
3. Econometrics provides tools for estimating causal effects using either observational (nonexperimental) data or data from real-world, imperfect experiments.
4. Econometrics also provides tools for predicting the value of a variable of interest using information in other, related variables.
5. Cross-sectional data are gathered by observing multiple entities at a single point in time; time series data are gathered by observing a single entity at multiple points in time; and panel data are gathered by observing multiple entities, each of which is observed at multiple points in time.

### Key Terms

randomized controlled experiment (48)  
control group (48)

treatment group (48)  
causal effect (48)

prediction (49)

forecast (49)

experimental data (49)

observational data (49)

cross-sectional data (50)

observation number (51)

time series data (51)

panel data (52)

longitudinal data (52)

### **MyLab Economics Can Help You Get a Better Grade**

**MyLab Economics** If your exam were tomorrow, would you be ready? For each chapter, **MyLab Economics** Practice Tests and Study Plan help you prepare for your exams. You can also find the Exercises and all Review the Concepts Questions available now in **MyLab Economics**. To see how it works, turn to the **MyLab Economics** spread on the inside front cover of this text and then go to [www.pearson.com/mylab/economics](http://www.pearson.com/mylab/economics).

For additional Empirical Exercises and Data Sets, log on to the Companion Website at [www.pearsonglobaleditions.com](http://www.pearsonglobaleditions.com).

## Review the Concepts

- 1.1** Describe a hypothetical ideal randomized controlled experiment to study the effect of six hours of reading on the improvement of the vocabulary of high school students. Suggest some impediments to implementing this experiment in practice.
- 1.2** Describe a hypothetical ideal randomized controlled experiment to study the effect of the consumption of alcohol on long-term memory loss. Suggest some impediments to implementing this experiment in practice.
- 1.3** You are asked to study the causal effect of hours spent on employee training (measured in hours per worker per week) in a manufacturing plant on the productivity of its workers (output per worker per hour). Describe:
  - a.** an ideal randomized controlled experiment to measure this causal effect;
  - b.** an observational cross-sectional data set with which you could study this effect;
  - c.** an observational time series data set for studying this effect; and
  - d.** an observational panel data set for studying this effect.

## Review of Probability

This chapter reviews the core ideas of the theory of probability that are needed to understand regression analysis and econometrics. We assume that you have taken an introductory course in probability and statistics. If your knowledge of probability is stale, you should refresh it by reading this chapter. If you feel confident with the material, you still should skim the chapter and the terms and concepts at the end to make sure you are familiar with the ideas and notation.

Most aspects of the world around us have an element of randomness. The theory of probability provides mathematical tools for quantifying and describing this randomness. Section 2.1 reviews probability distributions for a single random variable, and Section 2.2 covers the mathematical expectation, mean, and variance of a single random variable. Most of the interesting problems in economics involve more than one variable, and Section 2.3 introduces the basic elements of probability theory for two random variables. Section 2.4 discusses three special probability distributions that play a central role in statistics and econometrics: the normal, chi-squared, and  $F$  distributions.

The final two sections of this chapter focus on a specific source of randomness of central importance in econometrics: the randomness that arises by randomly drawing a sample of data from a larger population. For example, suppose you survey ten recent college graduates selected at random, record (or “observe”) their earnings, and compute the average earnings using these ten data points (or “observations”). Because you chose the sample at random, you could have chosen ten different graduates by pure random chance; had you done so, you would have observed ten different earnings, and you would have computed a different sample average. Because the average earnings vary from one randomly chosen sample to the next, the sample average is itself a random variable. Therefore, the sample average has a probability distribution, which is referred to as its sampling distribution because this distribution describes the different possible values of the sample average that would have occurred had a different sample been drawn.

Section 2.5 discusses random sampling and the sampling distribution of the sample average. This sampling distribution is, in general, complicated. When the sample size is sufficiently large, however, the sampling distribution of the sample average is approximately normal, a result known as the central limit theorem, which is discussed in Section 2.6.

## 2.1 Random Variables and Probability Distributions

### Probabilities, the Sample Space, and Random Variables

**Probabilities and outcomes.** The sex of the next new person you meet, your grade on an exam, and the number of times your wireless network connection fails while you are writing a term paper all have an element of chance or randomness. In each of these examples, there is something not yet known that is eventually revealed.

The mutually exclusive potential results of a random process are called the **outcomes**. For example, while writing your term paper, the wireless connection might never fail, it might fail once, it might fail twice, and so on. Only one of these outcomes will actually occur (the outcomes are mutually exclusive), and the outcomes need not be equally likely.

The **probability** of an outcome is the proportion of the time that the outcome occurs in the long run. If the probability of your wireless connection not failing while you are writing a term paper is 80%, then over the course of writing many term papers, you will complete 80% without a wireless connection failure.

**The sample space and events.** The set of all possible outcomes is called the **sample space**. An **event** is a subset of the sample space; that is, an event is a set of one or more outcomes. The event “my wireless connection will fail no more than once” is the set consisting of two outcomes: “no failures” and “one failure.”

**Random variables.** A random variable is a numerical summary of a random outcome. The number of times your wireless connection fails while you are writing a term paper is random and takes on a numerical value, so it is a random variable.

Some random variables are discrete and some are continuous. As their names suggest, a **discrete random variable** takes on only a discrete set of values, like 0, 1, 2, . . . , whereas a **continuous random variable** takes on a continuum of possible values.

### Probability Distribution of a Discrete Random Variable

**Probability distribution.** The **probability distribution** of a discrete random variable is the list of all possible values of the variable and the probability that each value will occur. These probabilities sum to 1.

For example, let  $M$  be the number of times your wireless network connection fails while you are writing a term paper. The probability distribution of the random variable  $M$  is the list of probabilities of all possible outcomes: The probability that  $M = 0$ , denoted  $\Pr(M = 0)$ , is the probability of no wireless connection failures;  $\Pr(M = 1)$  is the probability of a single connection failure; and so forth. An example of a probability distribution for  $M$  is given in the first row of Table 2.1. According to this distribution, the probability of no connection failures is 80%; the probability of one failure is 10%; and the probabilities of two, three, and four failures are,

**TABLE 2.1** Probability of Your Wireless Network Connection Failing  $M$  Times

	Outcome (number of failures)				
	0	1	2	3	4
Probability distribution	0.80	0.10	0.06	0.03	0.01
Cumulative probability distribution	0.80	0.90	0.96	0.99	1.00

respectively, 6%, 3%, and 1%. These probabilities sum to 100%. This probability distribution is plotted in Figure 2.1.

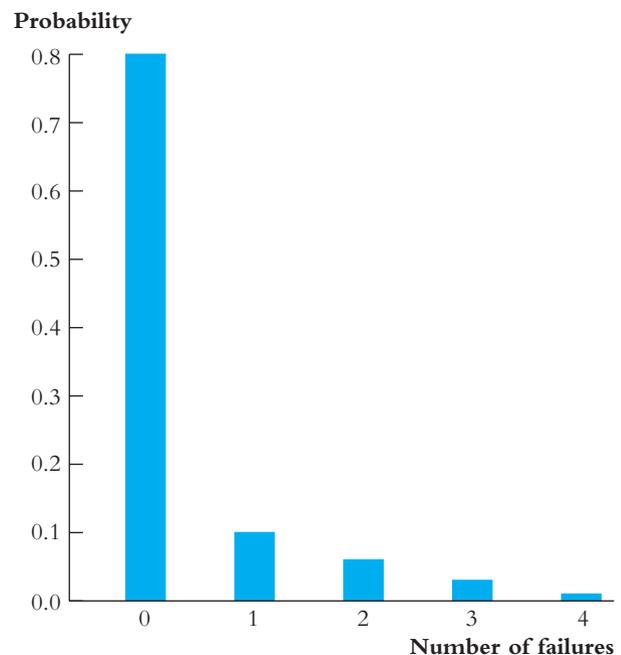
**Probabilities of events.** The probability of an event can be computed from the probability distribution. For example, the probability of the event of one or two failures is the sum of the probabilities of the constituent outcomes. That is,  $\Pr(M = 1 \text{ or } M = 2) = \Pr(M = 1) + \Pr(M = 2) = 0.10 + 0.06 = 0.16$ , or 16%.

**Cumulative probability distribution.** The **cumulative probability distribution** is the probability that the random variable is less than or equal to a particular value. The final row of Table 2.1 gives the cumulative probability distribution of the random variable  $M$ . For example, the probability of at most one connection failure,  $\Pr(M \leq 1)$ , is 90%, which is the sum of the probabilities of no failures (80%) and of one failure (10%).

A cumulative probability distribution is also referred to as a **cumulative distribution function**, a **c.d.f.**, or a **cumulative distribution**.

**FIGURE 2.1** Probability Distribution of the Number of Wireless Network Connection Failures

The height of each bar is the probability that the wireless connection fails the indicated number of times. The height of the first bar is 0.8, so the probability of 0 connection failures is 80%. The height of the second bar is 0.1, so the probability of 1 failure is 10%, and so forth for the other bars.



**The Bernoulli distribution.** An important special case of a discrete random variable is when the random variable is binary; that is, the outcome is 0 or 1. A binary random variable is called a **Bernoulli random variable** (in honor of the 17th-century Swiss mathematician and scientist Jacob Bernoulli), and its probability distribution is called the **Bernoulli distribution**.

For example, let  $G$  be the sex of the next new person you meet, where  $G = 0$  indicates that the person is male and  $G = 1$  indicates that the person is female. The outcomes of  $G$  and their probabilities thus are

$$G = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p, \end{cases} \quad (2.1)$$

where  $p$  is the probability of the next new person you meet being a woman. The probability distribution in Equation (2.1) is the Bernoulli distribution.

## Probability Distribution of a Continuous Random Variable

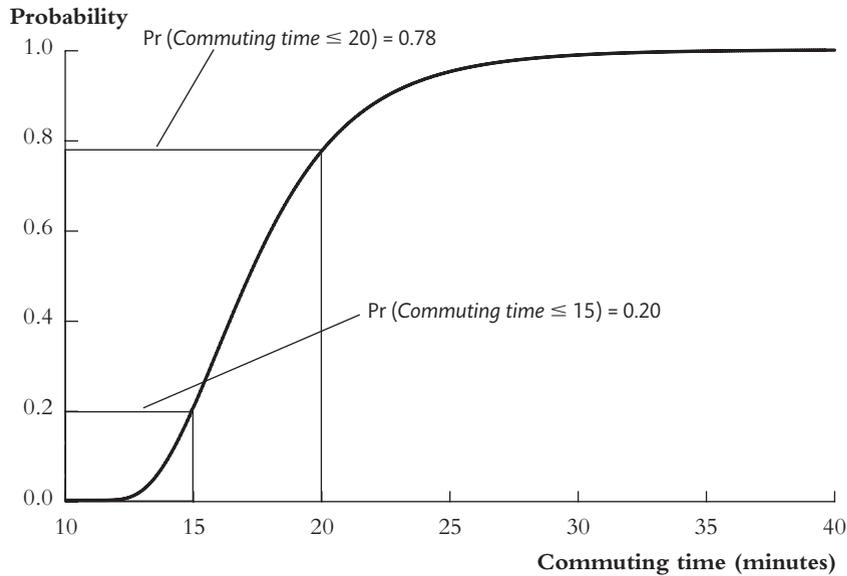
**Cumulative probability distribution.** The cumulative probability distribution for a continuous variable is defined just as it is for a discrete random variable. That is, the cumulative probability distribution of a continuous random variable is the probability that the random variable is less than or equal to a particular value.

For example, consider a student who drives from home to school. This student's commuting time can take on a continuum of values, and because it depends on random factors such as the weather and traffic conditions, it is natural to treat it as a continuous random variable. Figure 2.2a plots a hypothetical cumulative distribution of commuting times. For example, the probability that the commute takes less than 15 minutes is 20%, and the probability that it takes less than 20 minutes is 78%.

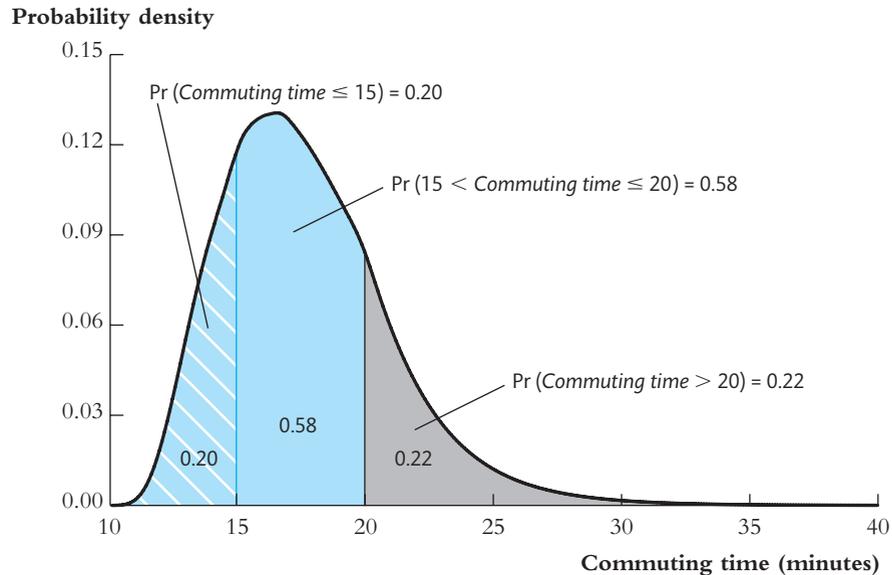
**Probability density function.** Because a continuous random variable can take on a continuum of possible values, the probability distribution used for discrete variables, which lists the probability of each possible value of the random variable, is not suitable for continuous variables. Instead, the probability is summarized by the **probability density function**. The area under the probability density function between any two points is the probability that the random variable falls between those two points. A probability density function is also called a **p.d.f.**, a **density function**, or simply a **density**.

Figure 2.2b plots the probability density function of commuting times corresponding to the cumulative distribution in Figure 2.2a. The probability that the commute takes between 15 and 20 minutes is given by the area under the p.d.f. between 15 minutes and 20 minutes, which is 0.58, or 58%. Equivalently, this probability can be seen on the cumulative distribution in Figure 2.2a as the difference between the probability that the commute is less than 20 minutes (78%) and the probability that it is less than 15 minutes (20%). Thus the probability density function and the cumulative probability distribution show the same information in different formats.

**FIGURE 2.2** Cumulative Probability Distribution and Probability Density Functions of Commuting Time



(a) Cumulative probability distribution function of commuting times



(b) Probability density function of commuting times

Figure 2.2a shows the cumulative probability distribution function (c.d.f.) of commuting times. The probability that a commuting time is less than 15 minutes is 0.20 (or 20%), and the probability that it is less than 20 minutes is 0.78 (78%). Figure 2.2b shows the probability density function (or p.d.f.) of commuting times. Probabilities are given by areas under the p.d.f. The probability that a commuting time is between 15 and 20 minutes is 0.58 (58%) and is given by the area under the curve between 15 and 20 minutes.

## 2.2 Expected Values, Mean, and Variance

### The Expected Value of a Random Variable

**Expected value.** The **expected value** of a random variable  $Y$ , denoted  $E(Y)$ , is the long-run average value of the random variable over many repeated trials or occurrences. The expected value of a discrete random variable is computed as a weighted average of the possible outcomes of that random variable, where the weights are the probabilities of that outcome. The expected value of  $Y$  is also called the **expectation** of  $Y$  or the **mean** of  $Y$  and is denoted  $\mu_Y$ .

For example, suppose you loan a friend \$100 at 10% interest. If the loan is repaid, you get \$110 (the principal of \$100 plus interest of \$10), but there is a risk of 1% that your friend will default and you will get nothing at all. Thus the amount you are repaid is a random variable that equals \$110 with probability 0.99 and equals \$0 with probability 0.01. Over many such loans, 99% of the time you would be paid back \$110, but 1% of the time you would get nothing, so on average you would be repaid  $\$110 \times 0.99 + \$0 \times 0.01 = \$108.90$ . Thus the expected value of your repayment is \$108.90.

As a second example, consider the number of wireless network connection failures  $M$  with the probability distribution given in Table 2.1. The expected value of  $M$ —that is, the mean of  $M$ —is the average number of failures over many term papers, weighted by the frequency with which a given number of failures occurs. Accordingly,

$$E(M) = 0 \times 0.80 + 1 \times 0.10 + 2 \times 0.06 + 3 \times 0.03 + 4 \times 0.01 = 0.35. \quad (2.2)$$

That is, the expected number of connection failures while writing a term paper is 0.35. Of course, the actual number of failures must always be an integer; it makes no sense to say that the wireless connection failed 0.35 times while writing a particular term paper! Rather, the calculation in Equation (2.2) means that the average number of failures over many such term papers is 0.35.

The formula for the expected value of a discrete random variable  $Y$  that can take on  $k$  different values is given in Key Concept 2.1. (Key Concept 2.1 uses summation notation, which is reviewed in Exercise 2.25.)

#### KEY CONCEPT

### Expected Value and the Mean

## 2.1

Suppose that the random variable  $Y$  takes on  $k$  possible values,  $y_1, \dots, y_k$ , where  $y_1$  denotes the first value,  $y_2$  denotes the second value, and so forth, and that the probability that  $Y$  takes on  $y_1$  is  $p_1$ , the probability that  $Y$  takes on  $y_2$  is  $p_2$ , and so forth. The expected value of  $Y$ , denoted  $E(Y)$ , is

$$E(Y) = y_1p_1 + y_2p_2 + \cdots + y_kp_k = \sum_{i=1}^k y_i p_i, \quad (2.3)$$

where the notation  $\sum_{i=1}^k y_i p_i$  means “the sum of  $y_i p_i$  for  $i$  running from 1 to  $k$ .” The expected value of  $Y$  is also called the mean of  $Y$  or the expectation of  $Y$  and is denoted  $\mu_Y$ .

**Expected value of a Bernoulli random variable.** An important special case of the general formula in Key Concept 2.1 is the mean of a Bernoulli random variable. Let  $G$  be the Bernoulli random variable with the probability distribution in Equation (2.1). The expected value of  $G$  is

$$E(G) = 0 \times (1 - p) + 1 \times p = p. \quad (2.4)$$

Thus the expected value of a Bernoulli random variable is  $p$ , the probability that it takes on the value 1.

**Expected value of a continuous random variable.** The expected value of a continuous random variable is also the probability-weighted average of the possible outcomes of the random variable. Because a continuous random variable can take on a continuum of possible values, the formal mathematical definition of its expectation involves calculus and its definition is given in Appendix 18.1.

## The Standard Deviation and Variance

The variance and standard deviation measure the dispersion or the “spread” of a probability distribution. The **variance** of a random variable  $Y$ , denoted  $\text{var}(Y)$ , is the expected value of the square of the deviation of  $Y$  from its mean:  $\text{var}(Y) = E[(Y - \mu_Y)^2]$ .

Because the variance involves the square of  $Y$ , the units of the variance are the units of the square of  $Y$ , which makes the variance awkward to interpret. It is therefore common to measure the spread by the **standard deviation**, which is the square root of the variance and is denoted  $\sigma_Y$ . The standard deviation has the same units as  $Y$ . These definitions are summarized in Key Concept 2.2.

For example, the variance of the number of connection failures  $M$  is the probability-weighted average of the squared difference between  $M$  and its mean, 0.35:

$$\begin{aligned} \text{var}(M) &= (0 - 0.35)^2 \times 0.80 + (1 - 0.35)^2 \times 0.10 + (2 - 0.35)^2 \times 0.06 \\ &\quad + (3 - 0.35)^2 \times 0.03 + (4 - 0.35)^2 \times 0.01 = 0.6475. \end{aligned} \quad (2.5)$$

The standard deviation of  $M$  is the square root of the variance, so  $\sigma_M = \sqrt{0.6475} \cong 0.80$ .

## Variance and Standard Deviation

### KEY CONCEPT

## 2.2

The variance of the discrete random variable  $Y$ , denoted  $\sigma_Y^2$ , is

$$\sigma_Y^2 = \text{var}(Y) = E[(Y - \mu_Y)^2] = \sum_{i=1}^k (y_i - \mu_Y)^2 p_i. \quad (2.6)$$

The standard deviation of  $Y$  is  $\sigma_Y$ , the square root of the variance. The units of the standard deviation are the same as the units of  $Y$ .

**Variance of a Bernoulli random variable.** The mean of the Bernoulli random variable  $G$  with the probability distribution in Equation (2.1) is  $\mu_G = p$  [Equation (2.4)], so its variance is

$$\text{var}(G) = \sigma_G^2 = (0 - p)^2 \times (1 - p) + (1 - p)^2 \times p = p(1 - p). \quad (2.7)$$

Thus the standard deviation of a Bernoulli random variable is  $\sigma_G = \sqrt{p(1 - p)}$ .

## Mean and Variance of a Linear Function of a Random Variable

This section discusses random variables (say,  $X$  and  $Y$ ) that are related by a linear function. For example, consider an income tax scheme under which a worker is taxed at a rate of 20% on his or her earnings and then given a (tax-free) grant of \$2000. Under this tax scheme, after-tax earnings  $Y$  are related to pre-tax earnings  $X$  by the equation

$$Y = 2000 + 0.8X. \quad (2.8)$$

That is, after-tax earnings  $Y$  is 80% of pre-tax earnings  $X$ , plus \$2000.

Suppose an individual's pre-tax earnings next year are a random variable with mean  $\mu_X$  and variance  $\sigma_X^2$ . Because pre-tax earnings are random, so are after-tax earnings. What are the mean and standard deviations of her after-tax earnings under this tax? After taxes, her earnings are 80% of the original pre-tax earnings, plus \$2000. Thus the expected value of her after-tax earnings is

$$E(Y) = \mu_Y = 2000 + 0.8\mu_X. \quad (2.9)$$

The variance of after-tax earnings is the expected value of  $(Y - \mu_Y)^2$ . Because  $Y = 2000 + 0.8X$ ,  $Y - \mu_Y = 2000 + 0.8X - (2000 + 0.8\mu_X) = 0.8(X - \mu_X)$ . Thus  $E[(Y - \mu_Y)^2] = E\{[0.8(X - \mu_X)]^2\} = 0.64E[(X - \mu_X)^2]$ . It follows that  $\text{var}(Y) = 0.64\text{var}(X)$ , so, taking the square root of the variance, the standard deviation of  $Y$  is

$$\sigma_Y = 0.8\sigma_X. \quad (2.10)$$

That is, the standard deviation of the distribution of her after-tax earnings is 80% of the standard deviation of the distribution of her pre-tax earnings.

This analysis can be generalized so that  $Y$  depends on  $X$  with an intercept  $a$  (instead of \$2000) and a slope  $b$  (instead of 0.8) so that

$$Y = a + bX. \quad (2.11)$$

Then the mean and variance of  $Y$  are

$$\mu_Y = a + b\mu_X \quad \text{and} \quad (2.12)$$

$$\sigma_Y^2 = b^2\sigma_X^2, \quad (2.13)$$

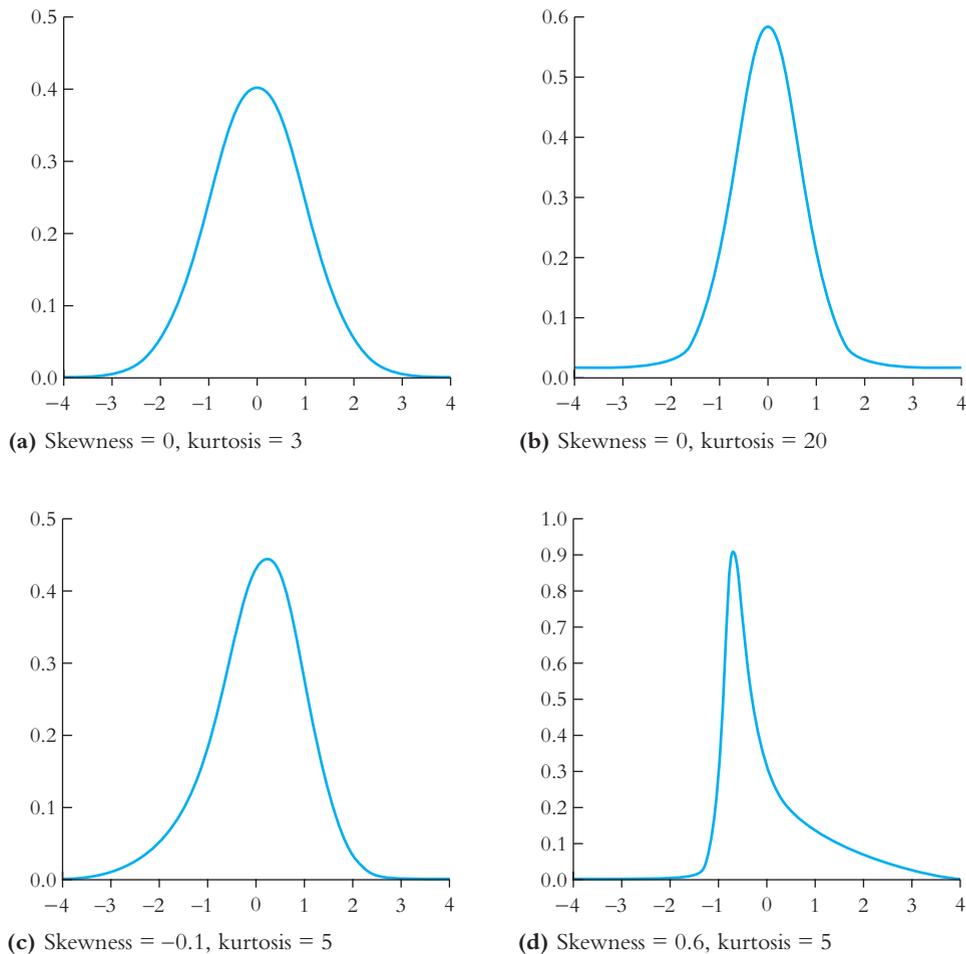
and the standard deviation of  $Y$  is  $\sigma_Y = b\sigma_X$ . The expressions in Equations (2.9) and (2.10) are applications of the more general formulas in Equations (2.12) and (2.13) with  $a = 2000$  and  $b = 0.8$ .

## Other Measures of the Shape of a Distribution

The mean and standard deviation measure two important features of a distribution: its center (the mean) and its spread (the standard deviation). This section discusses measures of two other features of a distribution: the skewness, which measures the lack of symmetry of a distribution, and the kurtosis, which measures how thick, or “heavy,” are its tails. The mean, variance, skewness, and kurtosis are all based on what are called the **moments of a distribution**.

**Skewness.** Figure 2.3 plots four distributions, two that are symmetric (Figures 2.3a and 2.3b) and two that are not (Figures 2.3c and 2.3d). Visually, the distribution in Figure 2.3d appears to deviate more from symmetry than does the distribution in

**FIGURE 2.3** Four Distributions with Different Skewness and Kurtosis



All of these distributions have a mean of 0 and a variance of 1. The distributions with skewness of 0 (a and b) are symmetric; the distributions with nonzero skewness (c and d) are not symmetric. The distributions with kurtosis exceeding 3 (b, c, and d) have heavy tails.

Figure 2.3c. The skewness of a distribution provides a mathematical way to describe how much a distribution deviates from symmetry.

The **skewness** of the distribution of a random variable  $Y$  is

$$\text{Skewness} = \frac{E[(Y - \mu_Y)^3]}{\sigma_Y^3}, \quad (2.14)$$

where  $\sigma_Y$  is the standard deviation of  $Y$ . For a symmetric distribution, a value of  $Y$  a given amount above its mean is just as likely as a value of  $Y$  the same amount below its mean. If so, then positive values of  $(Y - \mu_Y)^3$  will be offset on average (in expectation) by equally likely negative values. Thus, for a symmetric distribution,  $E(Y - \mu_Y)^3 = 0$ : The skewness of a symmetric distribution is 0. If a distribution is not symmetric, then a positive value of  $(Y - \mu_Y)^3$  generally is not offset on average by an equally likely negative value, so the skewness is nonzero for a distribution that is not symmetric. Dividing by  $\sigma_Y^3$  in the denominator of Equation (2.14) cancels the units of  $Y^3$  in the numerator, so the skewness is unit free; in other words, changing the units of  $Y$  does not change its skewness.

Below each of the four distributions in Figure 2.3 is its skewness. If a distribution has a long right tail, positive values of  $(Y - \mu_Y)^3$  are not fully offset by negative values, and the skewness is positive. If a distribution has a long left tail, its skewness is negative.

**Kurtosis.** The **kurtosis** of a distribution is a measure of how much mass is in its tails and therefore is a measure of how much of the variance of  $Y$  arises from extreme values. An extreme value of  $Y$  is called an **outlier**. The greater the kurtosis of a distribution, the more likely are outliers.

The kurtosis of the distribution of  $Y$  is

$$\text{Kurtosis} = \frac{E[(Y - \mu_Y)^4]}{\sigma_Y^4}. \quad (2.15)$$

If a distribution has a large amount of mass in its tails, then some extreme departures of  $Y$  from its mean are likely, and these departures will lead to large values, on average (in expectation), of  $(Y - \mu_Y)^4$ . Thus, for a distribution with a large amount of mass in its tails, the kurtosis will be large. Because  $(Y - \mu_Y)^4$  cannot be negative, the kurtosis cannot be negative.

The kurtosis of a normally distributed random variable is 3, so a random variable with kurtosis exceeding 3 has more mass in its tails than a normal random variable. A distribution with kurtosis exceeding 3 is called **leptokurtic** or, more simply, heavy-tailed. Like skewness, the kurtosis is unit free, so changing the units of  $Y$  does not change its kurtosis.

Below each of the four distributions in Figure 2.3 is its kurtosis. The distributions in Figures 2.3b–d are heavy-tailed.

**Moments.** The mean of  $Y$ ,  $E(Y)$ , is also called the first moment of  $Y$ , and the expected value of the square of  $Y$ ,  $E(Y^2)$ , is called the second moment of  $Y$ . In general, the

expected value of  $Y^r$  is called the  $r^{\text{th}}$  **moment** of the random variable  $Y$ . That is, the  $r^{\text{th}}$  moment of  $Y$  is  $E(Y^r)$ . The skewness is a function of the first, second, and third moments of  $Y$ , and the kurtosis is a function of the first through fourth moments of  $Y$ .

## Standardized Random Variables

A random variable can be transformed into a random variable with mean 0 and variance 1 by subtracting its mean and then dividing by its standard deviation, a process called standardization. Specifically, let  $Y$  have mean  $\mu_Y$  and variance  $\sigma_Y^2$ . Then the **standardized random variable** computed from  $Y$  is  $(Y - \mu_Y)/\sigma_Y$ . The mean of the standardized random variable is  $E(Y - \mu_Y)/\sigma_Y = (EY - \mu_Y)/\sigma_Y = 0$ , and its variance is  $\text{var}[(Y - \mu_Y)/\sigma_Y] = \text{var}(Y)/\sigma_Y^2 = 1$ . Standardized random variables do not have any units, such as dollars or meters, because the units of  $Y$  are canceled by dividing through by  $\sigma_Y$ , which also has the units of  $Y$ .

## 2.3 Two Random Variables

Most of the interesting questions in economics involve two or more variables. Are college graduates more likely to have a job than nongraduates? How does the distribution of income for women compare to that for men? These questions concern the distribution of two random variables, considered together (education and employment status in the first example, income and sex in the second). Answering such questions requires an understanding of the concepts of joint, marginal, and conditional probability distributions.

### Joint and Marginal Distributions

**Joint distribution.** The **joint probability distribution** of two discrete random variables, say  $X$  and  $Y$ , is the probability that the random variables simultaneously take on certain values, say  $x$  and  $y$ . The probabilities of all possible  $(x, y)$  combinations sum to 1. The joint probability distribution can be written as the function  $\text{Pr}(X = x, Y = y)$ .

For example, weather conditions—whether or not it is raining—affect the commuting time of the student commuter in Section 2.1. Let  $Y$  be a binary random variable that equals 1 if the commute is short (less than 20 minutes) and that equals 0 otherwise, and let  $X$  be a binary random variable that equals 0 if it is raining and 1 if not. Between these two random variables, there are four possible outcomes: it rains and the commute is long ( $X = 0, Y = 0$ ); rain and short commute ( $X = 0, Y = 1$ ); no rain and long commute ( $X = 1, Y = 0$ ); and no rain and short commute ( $X = 1, Y = 1$ ). The joint probability distribution is the frequency with which each of these four outcomes occurs over many repeated commutes.

An example of a joint distribution of these two variables is given in Table 2.2. According to this distribution, over many commutes, 15% of the days have rain and a long commute ( $X = 0, Y = 0$ ); that is, the probability of a long rainy commute is

**TABLE 2.2** Joint Distribution of Weather Conditions and Commuting Times

	Rain ( $X = 0$ )	No Rain ( $X = 1$ )	Total
Long commute ( $Y = 0$ )	0.15	0.07	0.22
Short commute ( $Y = 1$ )	0.15	0.63	0.78
Total	0.30	0.70	1.00

15%, or  $\Pr(X = 0, Y = 0) = 0.15$ . Also,  $\Pr(X = 0, Y = 1) = 0.15$ ,  $\Pr(X = 1, Y = 0) = 0.07$ , and  $\Pr(X = 1, Y = 1) = 0.63$ . These four possible outcomes are mutually exclusive and constitute the sample space, so the four probabilities sum to 1.

**Marginal probability distribution.** The **marginal probability distribution** of a random variable  $Y$  is just another name for its probability distribution. This term is used to distinguish the distribution of  $Y$  alone (the marginal distribution) from the joint distribution of  $Y$  and another random variable.

The marginal distribution of  $Y$  can be computed from the joint distribution of  $X$  and  $Y$  by adding up the probabilities of all possible outcomes for which  $Y$  takes on a specified value. If  $X$  can take on  $l$  different values  $x_1, \dots, x_l$ , then the marginal probability that  $Y$  takes on the value  $y$  is

$$\Pr(Y = y) = \sum_{i=1}^l \Pr(X = x_i, Y = y). \quad (2.16)$$

For example, in Table 2.2, the probability of a long rainy commute is 15%, and the probability of a long commute with no rain is 7%, so the probability of a long commute (rainy or not) is 22%. The marginal distribution of commuting times is given in the final column of Table 2.2. Similarly, the marginal probability that it will rain is 30%, as shown in the final row of Table 2.2.

## Conditional Distributions

**Conditional distribution.** The distribution of a random variable  $Y$  conditional on another random variable  $X$  taking on a specific value is called the **conditional distribution** of  $Y$  given  $X$ . The conditional probability that  $Y$  takes on the value  $y$  when  $X$  takes on the value  $x$  is written  $\Pr(Y = y | X = x)$ .

For example, what is the probability of a long commute ( $Y = 0$ ) if you know it is raining ( $X = 0$ )? From Table 2.2, the joint probability of a rainy short commute is 15%, and the joint probability of a rainy long commute is 15%, so if it is raining, a long commute and a short commute are equally likely. Thus the probability of a long commute ( $Y = 0$ ) conditional on it being rainy ( $X = 0$ ) is 50%, or  $\Pr(Y = 0 | X = 0) = 0.50$ . Equivalently, the marginal probability of rain is 30%; that is, over many commutes, it rains 30% of the time. Of this 30% of commutes, 50% of the time the commute is long ( $0.15/0.30$ ).

**TABLE 2.3** Joint and Conditional Distributions of Number of Wireless Connection Failures ( $M$ ) and Network Age ( $A$ )

<b>A. Joint Distribution</b>						
	$M = 0$	$M = 1$	$M = 2$	$M = 3$	$M = 4$	Total
Old network ( $A = 0$ )	0.35	0.065	0.05	0.025	0.01	0.50
New network ( $A = 1$ )	0.45	0.035	0.01	0.005	0.00	0.50
<b>Total</b>	0.80	0.10	0.06	0.03	0.01	1.00
<b>B. Conditional Distributions of <math>M</math> given <math>A</math></b>						
	$M = 0$	$M = 1$	$M = 2$	$M = 3$	$M = 4$	Total
$\Pr(M A = 0)$	0.70	0.13	0.10	0.05	0.02	1.00
$\Pr(M A = 1)$	0.90	0.07	0.02	0.01	0.00	1.00

In general, the conditional distribution of  $Y$  given  $X = x$  is

$$\Pr(Y = y|X = x) = \frac{\Pr(X = x, Y = y)}{\Pr(X = x)}. \quad (2.17)$$

For example, the conditional probability of a long commute given that it is rainy is  $\Pr(Y = 0|X = 0) = \Pr(X = 0, Y = 0)/\Pr(X = 0) = 0.15/0.30 = 0.50$ .

As a second example, consider a modification of the network connection failure example. Suppose that half the time you write your term paper in the school library, which has a new wireless network; otherwise, you write it in your room, which has an old wireless network. If we treat the location where you write the term paper as random, then the network age  $A$  ( $= 1$  if the network is new,  $= 0$  if it is old) is a random variable. Suppose the joint distribution of the random variables  $M$  and  $A$  is given in Part A of Table 2.3. Then the conditional distributions of connection failures given the age of the network are shown in Part B of the table. For example, the joint probability of  $M = 0$  and  $A = 0$  is 0.35; because half the time you use the old network, the conditional probability of no failures given that you use the old network is  $\Pr(M = 0|A = 0) = \Pr(M = 0, A = 0)/\Pr(A = 0) = 0.35/0.50 = 0.70$ , or 70%. In contrast, the conditional probability of no failures given that you use the new network is 90%. According to the conditional distributions in Part B of Table 2.3, the new network is less likely to fail than the old one; for example, the probability of three failures is 5% using the old network but 1% using the new network.

**Conditional expectation.** The **conditional expectation** of  $Y$  given  $X$ , also called the **conditional mean** of  $Y$  given  $X$ , is the mean of the conditional distribution of  $Y$  given  $X$ . That is, the conditional expectation is the expected value of  $Y$ , computed using the conditional distribution of  $Y$  given  $X$ . If  $Y$  takes on  $k$  values  $y_1, \dots, y_k$ , then the conditional mean of  $Y$  given  $X = x$  is

$$E(Y|X = x) = \sum_{i=1}^k y_i \Pr(Y = y_i|X = x). \quad (2.18)$$

For example, based on the conditional distributions in Table 2.3, the expected number of connection failures, given that the network is old, is  $E(M|A = 0) = 0 \times 0.70 + 1 \times 0.13 + 2 \times 0.10 + 3 \times 0.05 + 4 \times 0.02 = 0.56$ . The expected number of failures, given that the network is new, is  $E(M|A = 1) = 0.14$ , less than for the old network.

The conditional expectation of  $Y$  given  $X = x$  is just the mean value of  $Y$  when  $X = x$ . In the example of Table 2.3, the mean number of failures is 0.56 for the old network, so the conditional expectation of  $Y$  given that the network is old is 0.56. Similarly, for the new network, the mean number of failures is 0.14; that is, the conditional expectation of  $Y$  given that the network is new is 0.14.

**The law of iterated expectations.** The mean of  $Y$  is the weighted average of the conditional expectation of  $Y$  given  $X$ , weighted by the probability distribution of  $X$ . For example, the mean height of adults is the weighted average of the mean height of men and the mean height of women, weighted by the proportions of men and women. Stated mathematically, if  $X$  takes on the  $l$  values  $x_1, \dots, x_l$ , then

$$E(Y) = \sum_{i=1}^l E(Y|X = x_i) \Pr(X = x_i). \quad (2.19)$$

Equation (2.19) follows from Equations (2.18) and (2.17) (see Exercise 2.19).

Stated differently, the expectation of  $Y$  is the expectation of the conditional expectation of  $Y$  given  $X$ ,

$$E(Y) = E[E(Y|X)], \quad (2.20)$$

where the inner expectation on the right-hand side of Equation (2.20) is computed using the conditional distribution of  $Y$  given  $X$  and the outer expectation is computed using the marginal distribution of  $X$ . Equation (2.20) is known as the **law of iterated expectations**.

For example, the mean number of connection failures  $M$  is the weighted average of the conditional expectation of  $M$  given that it is old and the conditional expectation of  $M$  given that it is new, so  $E(M) = E(M|A = 0) \times \Pr(A = 0) + E(M|A = 1) \times \Pr(A = 1) = 0.56 \times 0.50 + 0.14 \times 0.50 = 0.35$ . This is the mean of the marginal distribution of  $M$ , as calculated in Equation (2.2).

The law of iterated expectations implies that if the conditional mean of  $Y$  given  $X$  is 0, then the mean of  $Y$  is 0. This is an immediate consequence of Equation (2.20): if  $E(Y|X) = 0$ , then  $E(Y) = E[E(Y|X)] = E[0] = 0$ . Said differently, if the mean of  $Y$  given  $X$  is 0, then it must be that the probability-weighted average of these conditional means is 0; that is, the mean of  $Y$  must be 0.

The law of iterated expectations also applies to expectations that are conditional on multiple random variables. For example, let  $X$ ,  $Y$ , and  $Z$  be random variables that are jointly distributed. Then the law of iterated expectations says that  $E(Y) = E[E(Y|X, Z)]$ , where  $E(Y|X, Z)$  is the conditional expectation of  $Y$

given both  $X$  and  $Z$ . For example, in the network connection illustration of Table 2.3, let  $P$  denote the number of people using the network; then  $E(M|A, P)$  is the expected number of failures for a network with age  $A$  that has  $P$  users. The expected number of failures overall,  $E(M)$ , is the weighted average of the expected number of failures for a network with age  $A$  and number of users  $P$ , weighted by the proportion of occurrences of both  $A$  and  $P$ .

Exercise 2.20 provides some additional properties of conditional expectations with multiple variables.

**Conditional variance.** The variance of  $Y$  conditional on  $X$  is the variance of the conditional distribution of  $Y$  given  $X$ . Stated mathematically, the **conditional variance** of  $Y$  given  $X$  is

$$\text{var}(Y|X = x) = \sum_{i=1}^k [y_i - E(Y|X = x)]^2 \Pr(Y = y_i|X = x). \quad (2.21)$$

For example, the conditional variance of the number of failures given that the network is old is  $\text{var}(M|A = 0) = (0 - 0.56)^2 \times 0.70 + (1 - 0.56)^2 \times 0.13 + (2 - 0.56)^2 \times 0.10 + (3 - 0.56)^2 \times 0.05 + (4 - 0.56)^2 \times 0.02 \cong 0.99$ . The standard deviation of the conditional distribution of  $M$  given that  $A = 0$  is thus  $\sqrt{0.99} = 0.99$ . The conditional variance of  $M$  given that  $A = 1$  is the variance of the distribution in the second row of Part B of Table 2.3, which is 0.22, so the standard deviation of  $M$  for the new network is  $\sqrt{0.22} = 0.47$ . For the conditional distributions in Table 2.3, the expected number of failures for the new network (0.14) is less than that for the old network (0.56), and the spread of the distribution of the number of failures, as measured by the conditional standard deviation, is smaller for the new network (0.47) than for the old (0.99).

**Bayes' rule.** **Bayes' rule** says that the conditional probability of  $Y$  given  $X$  is the conditional probability of  $X$  given  $Y$  times the relative marginal probabilities of  $Y$  and  $X$ :

$$\Pr(Y = y|X = x) = \frac{\Pr(X = x|Y = y)\Pr(Y = y)}{\Pr(X = x)} \text{ (Bayes' rule)}. \quad (2.22)$$

Equation (2.22) obtains from the definition of the conditional distribution in Equation (2.17), which implies that  $\Pr(X = x, Y = y) = \Pr(Y = y|X = x) \Pr(X = x)$  and that  $\Pr(X = x, Y = y) = \Pr(X = x|Y = y)\Pr(Y = y)$ ; equating the second parts of these two equalities and rearranging gives Bayes' rule.

Bayes' rule can be used to deduce conditional probabilities from the reverse conditional probability, with the help of marginal probabilities. For example, suppose you told your friend that you were dropped by the network three times last night while working on your term paper and your friend knows that half the time you work in the library and half the time you work in your room. Then your friend could deduce from Table 2.3 that the probability you worked in your room last night given three network failures is 83% (Exercise 2.28).

**The conditional mean is the minimum mean squared error prediction.** The conditional mean plays a central role in prediction; in fact it is, in a precise sense, the optimal prediction of  $Y$  given  $X = x$ .

A common formulation of the statistical prediction problem is to posit that the cost of making a prediction error increases with the square of that error. The motivation for this squared-error prediction loss is that small errors in prediction might not matter much, but large errors can be very costly in real-world applications. Stated mathematically, the prediction problem thus is: what is the function  $g(X)$  that minimizes the mean squared prediction error,  $E\{[Y - g(X)]^2\}$ ? The answer is the conditional mean  $E(Y|X)$ : Of all possible ways to use the information  $X$ , the conditional mean minimizes the mean squared prediction error. This result is proven in Appendix 2.2.

## Independence

Two random variables  $X$  and  $Y$  are **independently distributed**, or **independent**, if knowing the value of one of the variables provides no information about the other. Specifically,  $X$  and  $Y$  are independent if the conditional distribution of  $Y$  given  $X$  equals the marginal distribution of  $Y$ . That is,  $X$  and  $Y$  are independently distributed if, for all values of  $x$  and  $y$ ,

$$\Pr(Y = y|X = x) = \Pr(Y = y) \text{ (independence of } X \text{ and } Y). \quad (2.23)$$

Substituting Equation (2.23) into Equation (2.17) gives an alternative expression for independent random variables in terms of their joint distribution. If  $X$  and  $Y$  are independent, then

$$\Pr(X = x, Y = y) = \Pr(X = x)\Pr(Y = y). \quad (2.24)$$

That is, the joint distribution of two independent random variables is the product of their marginal distributions.

## Covariance and Correlation

**Covariance.** One measure of the extent to which two random variables move together is their covariance. The **covariance** between  $X$  and  $Y$  is the expected value  $E[(X - \mu_X)(Y - \mu_Y)]$ , where  $\mu_X$  is the mean of  $X$  and  $\mu_Y$  is the mean of  $Y$ . The covariance is denoted  $\text{cov}(X, Y)$  or  $\sigma_{XY}$ . If  $X$  can take on  $l$  values and  $Y$  can take on  $k$  values, then the covariance is given by the formula

$$\begin{aligned} \text{cov}(X, Y) &= \sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)] \\ &= \sum_{i=1}^k \sum_{j=1}^l (x_j - \mu_X)(y_i - \mu_Y)\Pr(X = x_j, Y = y_i). \end{aligned} \quad (2.25)$$

To interpret this formula, suppose that when  $X$  is greater than its mean (so that  $X - \mu_X$  is positive), then  $Y$  tends to be greater than its mean (so that  $Y - \mu_Y$  is

positive) and that when  $X$  is less than its mean (so that  $X - \mu_X < 0$ ), then  $Y$  tends to be less than its mean (so that  $Y - \mu_Y < 0$ ). In both cases, the product  $(X - \mu_X) \times (Y - \mu_Y)$  tends to be positive, so the covariance is positive. In contrast, if  $X$  and  $Y$  tend to move in opposite directions (so that  $X$  is large when  $Y$  is small, and vice versa), then the covariance is negative. Finally, if  $X$  and  $Y$  are independent, then the covariance is 0 (see Exercise 2.19).

**Correlation.** Because the covariance is the product of  $X$  and  $Y$ , deviated from their means, its units are, awkwardly, the units of  $X$  multiplied by the units of  $Y$ . This “units” problem can make numerical values of the covariance difficult to interpret.

The correlation is an alternative measure of dependence between  $X$  and  $Y$  that solves the “units” problem of the covariance. Specifically, the **correlation** between  $X$  and  $Y$  is the covariance between  $X$  and  $Y$  divided by their standard deviations:

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}. \quad (2.26)$$

Because the units of the numerator in Equation (2.26) are the same as those of the denominator, the units cancel, and the correlation is unit free. The random variables  $X$  and  $Y$  are said to be **uncorrelated** if  $\text{corr}(X, Y) = 0$ .

The correlation always is between  $-1$  and  $1$ ; that is, as proven in Appendix 2.1,

$$-1 \leq \text{corr}(X, Y) \leq 1 \quad (\text{correlation inequality}). \quad (2.27)$$

**Correlation and conditional mean.** If the conditional mean of  $Y$  does not depend on  $X$ , then  $Y$  and  $X$  are uncorrelated. That is,

$$\text{if } E(Y|X) = \mu_Y, \text{ then } \text{cov}(Y, X) = 0 \text{ and } \text{corr}(Y, X) = 0. \quad (2.28)$$

We now show this result. First, suppose  $Y$  and  $X$  have mean 0, so that  $\text{cov}(Y, X) = E[(Y - \mu_Y)(X - \mu_X)] = E(YX)$ . By the law of iterated expectations [Equation (2.20)],  $E(YX) = E[E(YX|X)] = E[E(Y|X)X] = 0$  because  $E(Y|X) = 0$ , so  $\text{cov}(Y, X) = 0$ . Equation (2.28) follows by substituting  $\text{cov}(Y, X) = 0$  into the definition of correlation in Equation (2.26). If  $Y$  and  $X$  do not have mean 0, subtract off their means, and then the preceding proof applies.

It is *not* necessarily true, however, that if  $X$  and  $Y$  are uncorrelated, then the conditional mean of  $Y$  given  $X$  does not depend on  $X$ . Said differently, it is possible for the conditional mean of  $Y$  to be a function of  $X$  but for  $Y$  and  $X$  nonetheless to be uncorrelated. An example is given in Exercise 2.23.

## The Mean and Variance of Sums of Random Variables

The mean of the sum of two random variables,  $X$  and  $Y$ , is the sum of their means:

$$E(X + Y) = E(X) + E(Y) = \mu_X + \mu_Y. \quad (2.29)$$

## The Distribution of Adulthood Earnings in the United Kingdom by Childhood Socioeconomic Circumstances

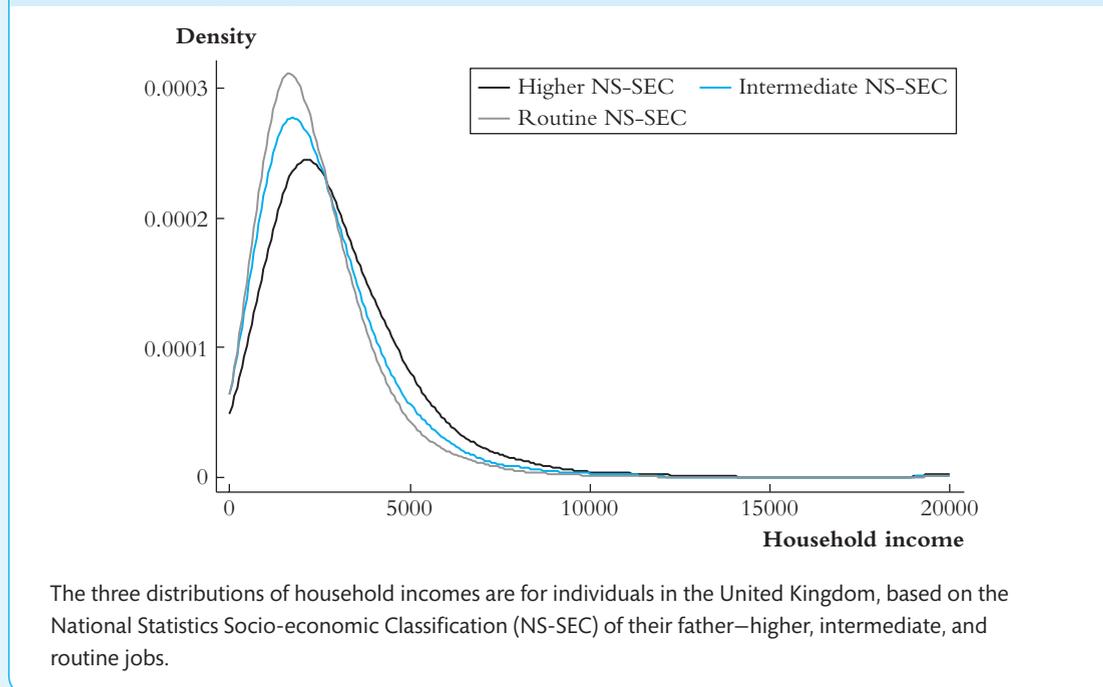
Politicians sometimes talk about how inequality in income arises as a result of differences in individual ability and effort. Are these politicians right? Or, in contrast, do childhood circumstances affect an individual's income during adulthood? For example, do children who grow up with fewer advantages go on to be part of households with lower average income?

One way to answer these questions is by considering how an individual's household income as

an adult varies according to their father's occupational type. While no two occupations are identical, researchers often group similar jobs into a given number of meaningful classes. One method of doing this, as seen in the United Kingdom's National Statistics Socio-economic Classification (NS-SEC),<sup>1</sup> is grouping jobs into a hierarchy of three classes: higher, intermediate, and routine.

Figure 2.4 illustrates these three conditional distributions of household income for individuals in

**FIGURE 2.4** Conditional Distributions of Household Income of U.K. individuals in 2009–2010, by Occupational Type of Father



<sup>1</sup>For further details refer to “The National Statistics Socio-economic classification (NS-SEC),” The Office for National Statistics, <https://www.ons.gov.uk/>, 2010.

**TABLE 2.4** Summaries of the Conditional Distribution of Monthly Household Income for Individuals in the United Kingdom Given NS-SEC of Father's Occupation

NS-SEC of Father's Job	Mean	Standard Deviation	Percentile			
			25%	50% (median)	75%	90%
(a) Higher	£3,149.27	£2,434.33	£1,663.33	£2,626.92	£3,973.74	£5,629.00
(b) Intermediate	2,692.01	2,187.53	1,362.44	2,237.56	3,382.00	4,881.99
(c) Routine	2,440.94	1,878.58	1,291.00	2,049.74	3,067.76	4,339.84

the United Kingdom in 2009 and 2010 according to the NS-SEC of their father's occupation in that individual's childhood.<sup>2</sup> The lower the classification of paternal occupation, the more concentrated in the lower end of the distribution is household income in adulthood.

The statistics for monthly household income for these individuals by NS-SEC classification are summarized in Table 2.4. For example, the mean income of individuals whose father's occupation is classified as routine, that is,  $E(\text{Income}|\text{Father's social class} = \text{routine})$ , was £2,440.94. This is over £700 less than that for individuals whose father's occupation is classified as higher, that is,  $E(\text{Income}|\text{Father's social class} = \text{higher})$ , which is £3,149.27. Furthermore, these differences are much greater at higher ends of the

distribution, with the difference in income between these groups being over £900 at the 75th percentile and almost £1,300 at the 90th percentile. The standard deviation of household income also increases with occupation classification, meaning that the spread of household income is also greater according to this measure.

This information is critical when examining the sort of claim discussed earlier. It appears that childhood circumstances may play some part in determining an individual's socioeconomic circumstances later in life. Can we say this for certain? Is there anything more to consider? These circumstances and others like a "gender gap" in earnings are an important aspect of the distribution of income. We revisit this topic in later chapters.

<sup>2</sup>Conditional distributions were estimated from data from the first wave of the United Kingdom's Understanding Society dataset (gathered during 2009 and 2010). More details are available at <https://www.understandingsociety.ac.uk/>. Individuals with missing observations are excluded.

## KEY CONCEPT

## 2.3

## Means, Variances, and Covariances of Sums of Random Variables

Let  $X$ ,  $Y$ , and  $V$  be random variables; let  $\mu_X$  and  $\sigma_X^2$  be the mean and variance of  $X$  and let  $\sigma_{XY}$  be the covariance between  $X$  and  $Y$  (and so forth for the other variables); and let  $a$ ,  $b$ , and  $c$  be constants. Equations (2.30) through (2.36) follow from the definitions of the mean, variance, and covariance:

$$E(a + bX + cY) = a + b\mu_X + c\mu_Y, \quad (2.30)$$

$$\text{var}(a + bY) = b^2\sigma_Y^2, \quad (2.31)$$

$$\text{var}(aX + bY) = a^2\sigma_X^2 + 2ab\sigma_{XY} + b^2\sigma_Y^2, \quad (2.32)$$

$$E(Y^2) = \sigma_Y^2 + \mu_Y^2, \quad (2.33)$$

$$\text{cov}(a + bX + cV, Y) = b\sigma_{XY} + c\sigma_{VY}, \quad (2.34)$$

$$E(XY) = \sigma_{XY} + \mu_X\mu_Y, \quad (2.35)$$

$$|\text{corr}(X, Y)| \leq 1 \text{ and } |\sigma_{XY}| \leq \sqrt{\sigma_X^2\sigma_Y^2} \text{ (correlation inequality)}. \quad (2.36)$$

The variance of the sum of  $X$  and  $Y$  is the sum of their variances plus two times their covariance:

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2\text{cov}(X, Y) = \sigma_X^2 + \sigma_Y^2 + 2\sigma_{XY}. \quad (2.37)$$

If  $X$  and  $Y$  are independent, then the covariance is 0, and the variance of their sum is the sum of their variances:

$$\begin{aligned} \text{var}(X + Y) &= \text{var}(X) + \text{var}(Y) = \sigma_X^2 + \sigma_Y^2 \\ &\text{(if } X \text{ and } Y \text{ are independent)}. \end{aligned} \quad (2.38)$$

Useful expressions for means, variances, and covariances involving weighted sums of random variables are collected in Key Concept 2.3. The results in Key Concept 2.3 are derived in Appendix 2.1.

## 2.4 The Normal, Chi-Squared, Student $t$ , and $F$ Distributions

The probability distributions most often encountered in econometrics are the normal, chi-squared, Student  $t$ , and  $F$  distributions.

### The Normal Distribution

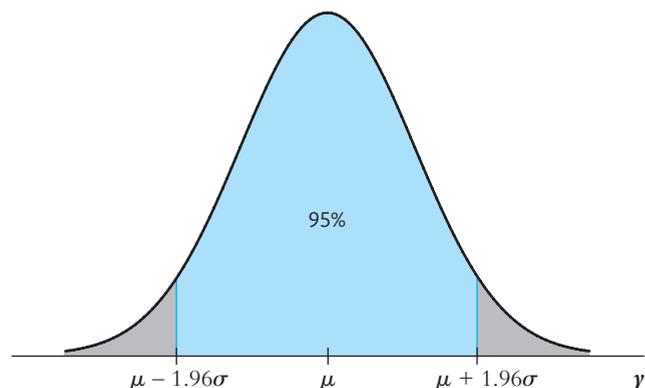
A continuous random variable with a **normal distribution** has the familiar bell-shaped probability density shown in Figure 2.5. The function defining the normal probability density is given in Appendix 18.1. As Figure 2.5 shows, the normal density with mean  $\mu$  and variance  $\sigma^2$  is symmetric around its mean and has 95% of its probability between  $\mu - 1.96\sigma$  and  $\mu + 1.96\sigma$ .

Some special notation and terminology have been developed for the normal distribution. The normal distribution with mean  $\mu$  and variance  $\sigma^2$  is expressed concisely as  $N(\mu, \sigma^2)$ . The **standard normal distribution** is the normal distribution with mean  $\mu = 0$  and variance  $\sigma^2 = 1$  and is denoted  $N(0, 1)$ . Random variables that have a  $N(0, 1)$  distribution are often denoted  $Z$ , and the standard normal cumulative distribution function is denoted by the Greek letter  $\Phi$ ; accordingly,  $\Pr(Z \leq c) = \Phi(c)$ , where  $c$  is a constant. Values of the standard normal cumulative distribution function are tabulated in Appendix Table 1.

To look up probabilities for a normal variable with a general mean and variance, we must first standardize the variable. For example, suppose  $Y$  is distributed  $N(1, 4)$ —that is,  $Y$  is normally distributed with a mean of 1 and a variance of 4. What is the probability that  $Y \leq 2$ —that is, what is the shaded area in Figure 2.6a? The standardized version of  $Y$  is  $Y$  minus its mean, divided by its standard deviation; that is,  $(Y - 1)/\sqrt{4} = \frac{1}{2}(Y - 1)$ . Accordingly, the random variable  $\frac{1}{2}(Y - 1)$  is normally distributed with mean 0 and variance 1 (see Exercise 2.8); it has the standard normal

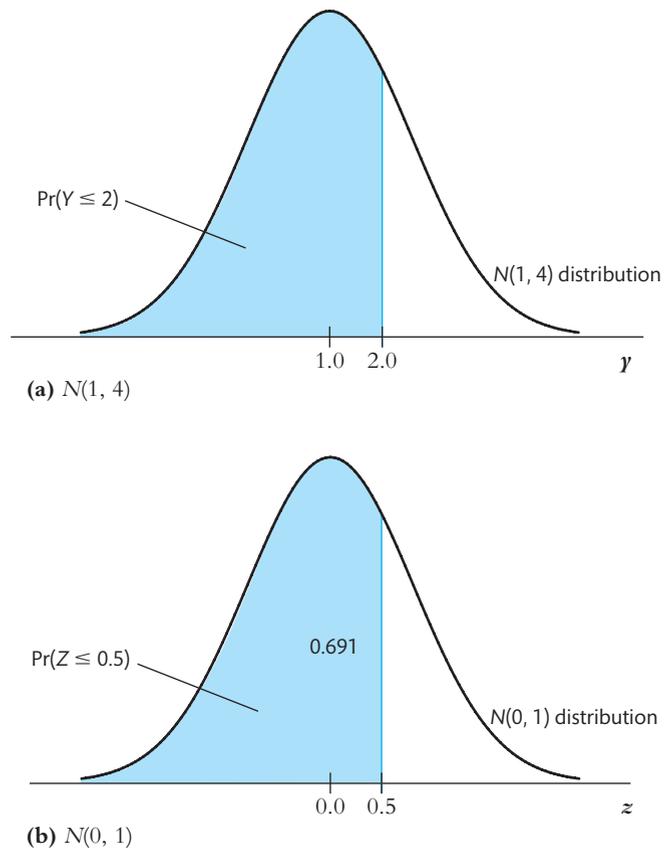
**FIGURE 2.5** The Normal Probability Density

The normal probability density function with mean  $\mu$  and variance  $\sigma^2$  is a bell-shaped curve, centered at  $\mu$ . The area under the normal p.d.f. between  $\mu - 1.96\sigma$  and  $\mu + 1.96\sigma$  is 0.95. The normal distribution is denoted  $N(\mu, \sigma^2)$ .



**FIGURE 2.6** Calculating the Probability That  $Y \leq 2$  When  $Y$  Is Distributed  $N(1, 4)$ 

To calculate  $\Pr(Y \leq 2)$ , standardize  $Y$ , then use the standard normal distribution table.  $Y$  is standardized by subtracting its mean ( $\mu = 1$ ) and dividing by its standard deviation ( $\sigma = 2$ ). The probability that  $Y \leq 2$  is shown in Figure 2.6a, and the corresponding probability after standardizing  $Y$  is shown in Figure 2.6b. Because the standardized random variable,  $(Y - 1)/2$ , is a standard normal ( $Z$ ) random variable,  $\Pr(Y \leq 2) = \Pr\left(\frac{Y-1}{2} \leq \frac{2-1}{2}\right) = \Pr(Z \leq 0.5)$ . From Appendix Table 1,  $\Pr(Z \leq 0.5) = \Phi(0.5) = 0.691$ .

**KEY CONCEPT**

## 2.4

**Computing Probabilities and Involving Normal Random Variables**

Suppose  $Y$  is normally distributed with mean  $\mu$  and variance  $\sigma^2$ ; in other words,  $Y$  is distributed  $N(\mu, \sigma^2)$ . Then  $Y$  is standardized by subtracting its mean and dividing by its standard deviation, that is, by computing  $Z = (Y - \mu)/\sigma$ .

Let  $c_1$  and  $c_2$  denote two numbers with  $c_1 < c_2$ , and let  $d_1 = (c_1 - \mu)/\sigma$  and  $d_2 = (c_2 - \mu)/\sigma$ . Then

$$\Pr(Y \leq c_2) = \Pr(Z \leq d_2) = \Phi(d_2), \quad (2.39)$$

$$\Pr(Y \geq c_1) = \Pr(Z \geq d_1) = 1 - \Phi(d_1), \quad (2.40)$$

$$\Pr(c_1 \leq Y \leq c_2) = \Pr(d_1 \leq Z \leq d_2) = \Phi(d_2) - \Phi(d_1). \quad (2.41)$$

The normal cumulative distribution function  $\Phi$  is tabulated in Appendix Table 1.

distribution shown in Figure 2.6b. Now  $Y \leq 2$  is equivalent to  $\frac{1}{2}(Y - 1) \leq \frac{1}{2}(2 - 1)$ ; that is,  $\frac{1}{2}(Y - 1) \leq \frac{1}{2}$ . Thus

$$\Pr(Y \leq 2) = \Pr\left[\frac{1}{2}(Y - 1) \leq \frac{1}{2}\right] = \Pr(Z \leq \frac{1}{2}) = \Phi(0.5) = 0.691, \quad (2.42)$$

where the value 0.691 is taken from Appendix Table 1.

The same approach can be used to compute the probability that a normally distributed random variable exceeds (or is less than) some value or that it falls in a certain range. These steps are discussed in Key Concept 2.4. The box “The Unpegging of the Swiss Franc” presents an unusual application of the cumulative normal distribution.

The normal distribution is symmetric, so its skewness is 0. The kurtosis of the normal distribution is 3.

**The multivariate normal distribution.** The normal distribution can be generalized to describe the joint distribution of a set of random variables. In this case, the distribution is called the **multivariate normal distribution** or, if only two variables are being considered, the **bivariate normal distribution**. The formula for the bivariate normal p.d.f. is given in Appendix 18.1, and the formula for the general multivariate normal p.d.f. is given in Appendix 19.2.

The multivariate normal distribution has four important properties. If  $X$  and  $Y$  have a bivariate normal distribution with covariance  $\sigma_{XY}$  and if  $a$  and  $b$  are two constants, then  $aX + bY$  has the normal distribution:

$$aX + bY \text{ is distributed } N(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\sigma_{XY}) \\ (X, Y \text{ bivariate normal}). \quad (2.43)$$

## The Unpegging of the Swiss Franc

**O**n Thursday, January 15, 2015, the value of the euro fell by 17.472% from 1.201 to 0.991 against the Swiss franc. This was a huge shift, illustrated in the downward spike in Figure 2.7, given that the previous year had not seen a day’s movement greater than 0.544%. If you had woken up as a statistical analyst for a financial company on that Thursday morning, how might you have estimated the probability of this happening that day?

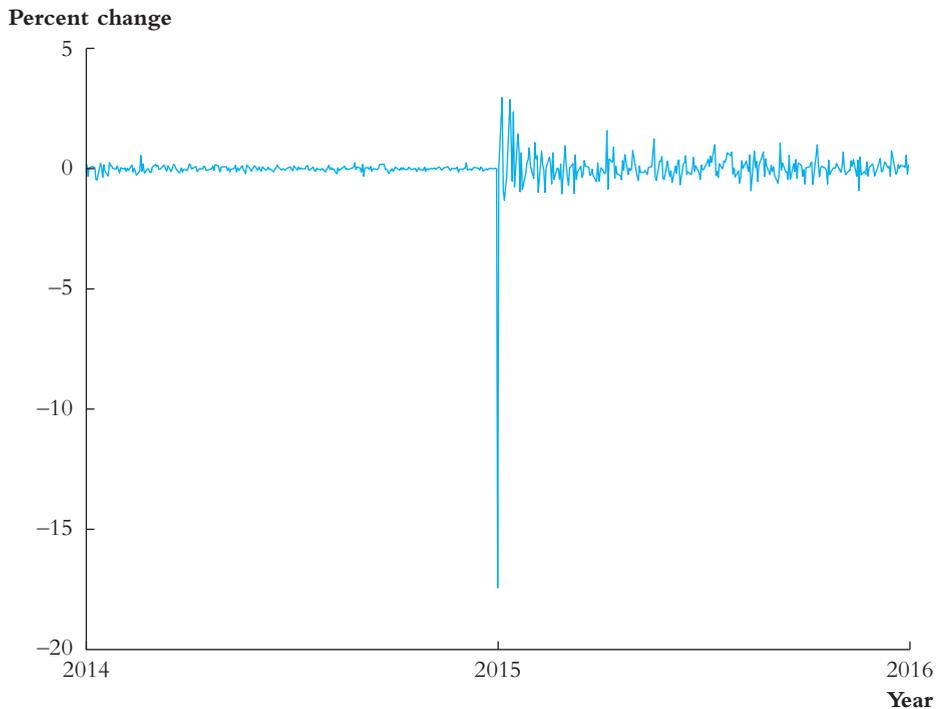
If you had assumed the data was normally distributed, you would have required an estimate of the standard deviation of daily percentage change in the euro/Swiss franc exchange rates. Using Datastream

data<sup>1</sup> for the year to January 14, 2015, you can estimate this as 0.112%.

What was the probability of a drop of 17.472%? We can first calculate the number of standard deviations that describes a change of this magnitude as  $\frac{17.472}{0.112} = 156$ . If the daily percentage changes are normally distributed, then the estimate of the probability of a fall at least as big as 156 standard deviations corresponds to an inconceivably small number— $8.175 \times 10^{-5288}$ , which is derived using Equation (2.39).

<sup>1</sup>Datastream, maintained by Thomson Reuters, is a global financial and macroeconomic data platform that acts as a repository of financial and economic data.

*continued on next page*

**FIGURE 2.7** Daily Percentage Change in the Euro/Swiss Franc Exchange Rate

The day-on-day percentage change in the value of the euro in Swiss francs for a year before and a year after the unpegging of the Swiss franc on January 15, 2015.

So was the probability of a fall at least this large really so small? Well, no. The error here is to not investigate the nature of our data further, and to fail to understand the actual process that determined the value of the currency. The Swiss franc had in fact been kept within very small bounds due to the actions of the country's central bank in setting a so-called "peg" for the currency. In the previous twelve months, this had been within the range of 1.2008 and 1.236 Swiss francs per euro. In fact, the introduction of this peg over three years earlier had caused an appreciation of the euro against the Swiss franc of over 20 standard deviations (again, assuming a normal distribution derived from previous daily changes!).<sup>2</sup>

It was the introduction of the peg that had caused such little volatility in—or such a low standard deviation of—the value of the currency. Once this peg was removed, as happened on that particular Thursday, the value of the currency was able to float and vary according to market factors. Investors responded to the removal of the peg by bidding down the value of the euro against the franc substantially.

It is not only the removal of a currency peg in this way that can cause extreme fluctuations. The result

<sup>2</sup>See the article published in Reuters, "Charts of the Day, Swiss Franc Edition," by Felix Salmon, September 6, 2011.

of the 2016 “Brexit” referendum in the United Kingdom—an event that, while seen as unlikely, was at least partly foreseeable—led to an appreciation in the value of the euro against British pound sterling on June 24, 2016, of 6.17%. This is equivalent to 9.80 standard deviations (based on data from the previous year), or an event with an apparent probability of  $5.629 \times 10^{-23}$ . While it may seem substantially more likely to occur, the probability of such an event actually taking place is less than once every 1,000,000,000,000,000,000 years (a total of 18 zeros)!<sup>3</sup> Again, it seems unlikely that this

is an accurate characterization of the probability of such an event occurring.

Clearly, it is dangerous to assume that data is normally distributed or that recent observations of a variable will provide a useful prediction of the range of future values. Indeed, it is partly for this reason that advertisements for financial products in the United Kingdom must carry a disclaimer that “past performance is not a guide to future performance.”

<sup>3</sup>This is based on the assumption of 260 trading days per year.

More generally, if  $n$  random variables have a multivariate normal distribution, then any linear combination of these variables (such as their sum) is normally distributed.

Second, if a set of variables has a multivariate normal distribution, then the marginal distribution of each of the variables is normal [this follows from Equation (2.43) by setting  $a = 1$  and  $b = 0$ ].

Third, if variables with a multivariate normal distribution have covariances that equal 0, then the variables are independent. Thus, if  $X$  and  $Y$  have a bivariate normal distribution and  $\sigma_{XY} = 0$ , then  $X$  and  $Y$  are independent (this is shown in Appendix 18.1). In Section 2.3, it was shown that if  $X$  and  $Y$  are independent, then, regardless of their joint distribution,  $\sigma_{XY} = 0$ . If  $X$  and  $Y$  are jointly normally distributed, then the converse is also true. This result—that 0 covariance implies independence—is a special property of the multivariate normal distribution that is not true in general.

Fourth, if  $X$  and  $Y$  have a bivariate normal distribution, then the conditional expectation of  $Y$  given  $X$  is linear in  $X$ ; that is,  $E(Y|X = x) = a + bx$ , where  $a$  and  $b$  are constants (Exercise 18.11). Joint normality implies linearity of conditional expectations, but linearity of conditional expectations does not imply joint normality.