

SIXTH EDITION



Introductory
Econometrics

A Modern Approach

JEFFREY M. WOOLDRIDGE

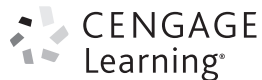
Introductory Econometrics

A MODERN APPROACH

SIXTH EDITION

Jeffrey M. Wooldridge

Michigan State University



Australia • Brazil • Mexico • Singapore • United Kingdom • United States

This is an electronic version of the print textbook. Due to electronic rights restrictions, some third party content may be suppressed. Editorial review has deemed that any suppressed content does not materially affect the overall learning experience. The publisher reserves the right to remove content from this title at any time if subsequent rights restrictions require it. For valuable information on pricing, previous editions, changes to current editions, and alternate formats, please visit www.cengage.com/highered to search by ISBN#, author, title, or keyword for materials in your areas of interest.

Important Notice: Media content referenced within the product description or the product text may not be available in the eBook version.

Introductory Econometrics, 6e
Jeffrey M. Wooldridge

Vice President, General Manager, Social
Science & Qualitative Business: Erin Joyner

Product Director: Mike Worls

Associate Product Manager: Tara Singer

Content Developer: Chris Rader

Marketing Director: Kristen Hurd

Marketing Manager: Katie Jergens

Marketing Coordinator: Chris Walz

Art and Cover Direction, Production
Management, and Composition: Lumina
Datamatics, Inc.

Intellectual Property Analyst: Jennifer
Nonenmacher

Project Manager: Sarah Shainwald

Manufacturing Planner: Kevin Kluck

Cover Image: ©kentoh/Shutterstock

Unless otherwise noted, all items
© Cengage Learning

© 2016, 2013 Cengage Learning

WCN: 02-200-203

ALL RIGHTS RESERVED. No part of this work covered by the copyright herein may be reproduced, transmitted, stored, or used in any form or by any means graphic, electronic, or mechanical, including but not limited to photocopying, recording, scanning, digitizing, taping, Web distribution, information networks, or information storage and retrieval systems, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without the prior written permission of the publisher.

For product information and technology assistance, contact us at **Cengage Learning Customer & Sales Support, 1-800-354-9706**

For permission to use material from this text or product, submit all requests online at www.cengage.com/permissions
Further permissions questions can be emailed to permissionrequest@cengage.com

Library of Congress Control Number: 2015944828

Student Edition:

ISBN: 978-1-305-27010-7

Cengage Learning

20 Channel Center Street

Boston, MA 02210

USA

Cengage Learning is a leading provider of customized learning solutions with employees residing in nearly 40 different countries and sales in more than 125 countries around the world. Find your local representative at **www.cengage.com**.

Cengage Learning products are represented in Canada by Nelson Education, Ltd.

To learn more about Cengage Learning Solutions, visit

www.cengage.com

Purchase any of our products at your local college store or at our preferred online store **www.cengagebrain.com**

Brief Contents

Chapter 1	The Nature of Econometrics and Economic Data	1
PART 1: Regression Analysis with Cross-Sectional Data		19
Chapter 2	The Simple Regression Model	20
Chapter 3	Multiple Regression Analysis: Estimation	60
Chapter 4	Multiple Regression Analysis: Inference	105
Chapter 5	Multiple Regression Analysis: OLS Asymptotics	149
Chapter 6	Multiple Regression Analysis: Further Issues	166
Chapter 7	Multiple Regression Analysis with Qualitative Information: Binary (or Dummy) Variables	205
Chapter 8	Heteroskedasticity	243
Chapter 9	More on Specification and Data Issues	274
PART 2: Regression Analysis with Time Series Data		311
Chapter 10	Basic Regression Analysis with Time Series Data	312
Chapter 11	Further Issues in Using OLS with Time Series Data	344
Chapter 12	Serial Correlation and Heteroskedasticity in Time Series Regressions	372
PART 3: Advanced Topics		401
Chapter 13	Pooling Cross Sections Across Time: Simple Panel Data Methods	402
Chapter 14	Advanced Panel Data Methods	434
Chapter 15	Instrumental Variables Estimation and Two Stage Least Squares	461
Chapter 16	Simultaneous Equations Models	499
Chapter 17	Limited Dependent Variable Models and Sample Selection Corrections	524
Chapter 18	Advanced Time Series Topics	568
Chapter 19	Carrying Out an Empirical Project	605
APPENDICES		
Appendix A	Basic Mathematical Tools	628
Appendix B	Fundamentals of Probability	645
Appendix C	Fundamentals of Mathematical Statistics	674
Appendix D	Summary of Matrix Algebra	709
Appendix E	The Linear Regression Model in Matrix Form	720
Appendix F	Answers to Chapter Questions	734
Appendix G	Statistical Tables	743
References		750
Glossary		756
Index		771

Contents

Preface xii
About the Author xxi

CHAPTER 1 The Nature of Econometrics and Economic Data 1

1-1 What Is Econometrics? 1
1-2 Steps in Empirical Economic Analysis 2
1-3 The Structure of Economic Data 5
 1-3a *Cross-Sectional Data* 5
 1-3b *Time Series Data* 7
 1-3c *Pooled Cross Sections* 8
 1-3d *Panel or Longitudinal Data* 9
 1-3e *A Comment on Data Structures* 10
1-4 Causality and the Notion of *Ceteris Paribus* in Econometric Analysis 10
Summary 14
Key Terms 14
Problems 15
Computer Exercises 15

PART 1

Regression Analysis with Cross-Sectional Data 19

CHAPTER 2 The Simple Regression Model 20

2-1 Definition of the Simple Regression Model 20
2-2 Deriving the Ordinary Least Squares Estimates 24
 2-2a *A Note on Terminology* 31
2-3 Properties of OLS on Any Sample of Data 32
 2-3a *Fitted Values and Residuals* 32
 2-3b *Algebraic Properties of OLS Statistics* 32
 2-3c *Goodness-of-Fit* 35

2-4 Units of Measurement and Functional Form 36
 2-4a *The Effects of Changing Units of Measurement on OLS Statistics* 36
 2-4b *Incorporating Nonlinearities in Simple Regression* 37
 2-4c *The Meaning of “Linear” Regression* 40
2-5 Expected Values and Variances of the OLS Estimators 40
 2-5a *Unbiasedness of OLS* 40
 2-5b *Variances of the OLS Estimators* 45
 2-5c *Estimating the Error Variance* 48
2-6 Regression through the Origin and Regression on a Constant 50
Summary 51
Key Terms 52
Problems 53
Computer Exercises 56
Appendix 2A 59

CHAPTER 3 Multiple Regression Analysis: Estimation 60

3-1 Motivation for Multiple Regression 61
 3-1a *The Model with Two Independent Variables* 61
 3-1b *The Model with k Independent Variables* 63
3-2 Mechanics and Interpretation of Ordinary Least Squares 64
 3-2a *Obtaining the OLS Estimates* 64
 3-2b *Interpreting the OLS Regression Equation* 65
 3-2c *On the Meaning of “Holding Other Factors Fixed” in Multiple Regression* 67
 3-2d *Changing More Than One Independent Variable Simultaneously* 68
 3-2e *OLS Fitted Values and Residuals* 68
 3-2f *A “Partialling Out” Interpretation of Multiple Regression* 69

3-2g	<i>Comparison of Simple and Multiple Regression Estimates</i>	69
3-2h	<i>Goodness-of-Fit</i>	70
3-2i	<i>Regression through the Origin</i>	73
3-3	The Expected Value of the OLS Estimators	73
3-3a	<i>Including Irrelevant Variables in a Regression Model</i>	77
3-3b	<i>Omitted Variable Bias: The Simple Case</i>	78
3-3c	<i>Omitted Variable Bias: More General Cases</i>	81
3-4	The Variance of the OLS Estimators	81
3-4a	<i>The Components of the OLS Variances. Multicollinearity</i>	83
3-4b	<i>Variances in Misspecified Models</i>	86
3-4c	<i>Estimating σ^2 Standard Errors of the OLS Estimators</i>	87
3-5	Efficiency of OLS: The Gauss-Markov Theorem	89
3-6	Some Comments on the Language of Multiple Regression Analysis	90
	Summary	91
	Key Terms	93
	Problems	93
	Computer Exercises	97
	Appendix 3A	101

CHAPTER 4 Multiple Regression Analysis: Inference 105

4-1	Sampling Distributions of the OLS Estimators	105
4-2	Testing Hypotheses about a Single Population Parameter: The t Test	108
4-2a	<i>Testing against One-Sided Alternatives</i>	110
4-2b	<i>Two-Sided Alternatives</i>	114
4-2c	<i>Testing Other Hypotheses about β_j</i>	116
4-2d	<i>Computing p-Values for t Tests</i>	118
4-2e	<i>A Reminder on the Language of Classical Hypothesis Testing</i>	120
4-2f	<i>Economic, or Practical, versus Statistical Significance</i>	120
4-3	Confidence Intervals	122
4-4	Testing Hypotheses about a Single Linear Combination of the Parameters	124
4-5	Testing Multiple Linear Restrictions: The F Test	127
4-5a	<i>Testing Exclusion Restrictions</i>	127
4-5b	<i>Relationship between F and t Statistics</i>	132
4-5c	<i>The R-Squared Form of the F Statistic</i>	133
4-5d	<i>Computing p-Values for F Tests</i>	134
4-5e	<i>The F Statistic for Overall Significance of a Regression</i>	135
4-5f	<i>Testing General Linear Restrictions</i>	136

4-6	Reporting Regression Results	137
	Summary	139
	Key Terms	140
	Problems	141
	Computer Exercises	146

CHAPTER 5 Multiple Regression Analysis: OLS Asymptotics 149

5-1	Consistency	150
5-1a	<i>Deriving the Inconsistency in OLS</i>	153
5-2	Asymptotic Normality and Large Sample Inference	154
5-2a	<i>Other Large Sample Tests: The Lagrange Multiplier Statistic</i>	158
5-3	Asymptotic Efficiency of OLS	161
	Summary	162
	Key Terms	162
	Problems	162
	Computer Exercises	163
	Appendix 5A	165

CHAPTER 6 Multiple Regression Analysis: Further Issues 166

6-1	Effects of Data Scaling on OLS Statistics	166
6-1a	<i>Beta Coefficients</i>	169
6-2	More on Functional Form	171
6-2a	<i>More on Using Logarithmic Functional Forms</i>	171
6-2b	<i>Models with Quadratics</i>	173
6-2c	<i>Models with Interaction Terms</i>	177
6-2d	<i>Computing Average Partial Effects</i>	179
6-3	More on Goodness-of-Fit and Selection of Regressors	180
6-3a	<i>Adjusted R-Squared</i>	181
6-3b	<i>Using Adjusted R-Squared to Choose between Nonnested Models</i>	182
6-3c	<i>Controlling for Too Many Factors in Regression Analysis</i>	184
6-3d	<i>Adding Regressors to Reduce the Error Variance</i>	185
6-4	Prediction and Residual Analysis	186
6-4a	<i>Confidence Intervals for Predictions</i>	186
6-4b	<i>Residual Analysis</i>	190
6-4c	<i>Predicting y When $\log(y)$ Is the Dependent Variable</i>	190
6-4d	<i>Predicting y When the Dependent Variable Is $\log(y)$</i>	192

Summary 194
 Key Terms 196
 Problems 196
 Computer Exercises 199
 Appendix 6A 203

CHAPTER 7 Multiple Regression Analysis with Qualitative Information: Binary (or Dummy) Variables 205

7-1 Describing Qualitative Information 205
 7-2 A Single Dummy Independent Variable 206
 7-2a *Interpreting Coefficients on Dummy Explanatory Variables When the Dependent Variable Is $\log(y)$* 211
 7-3 Using Dummy Variables for Multiple Categories 212
 7-3a *Incorporating Ordinal Information by Using Dummy Variables* 214
 7-4 Interactions Involving Dummy Variables 217
 7-4a *Interactions among Dummy Variables* 217
 7-4b *Allowing for Different Slopes* 218
 7-4c *Testing for Differences in Regression Functions across Groups* 221
 7-5 A Binary Dependent Variable: The Linear Probability Model 224
 7-6 More on Policy Analysis and Program Evaluation 229
 7-7 Interpreting Regression Results with Discrete Dependent Variables 231
 Summary 232
 Key Terms 233
 Problems 233
 Computer Exercises 237

CHAPTER 8 Heteroskedasticity 243

8-1 Consequences of Heteroskedasticity for OLS 243
 8-2 Heteroskedasticity-Robust Inference after OLS Estimation 244
 8-2a *Computing Heteroskedasticity-Robust LM Tests* 248
 8-3 Testing for Heteroskedasticity 250
 8-3a *The White Test for Heteroskedasticity* 252
 8-4 Weighted Least Squares Estimation 254
 8-4a *The Heteroskedasticity Is Known up to a Multiplicative Constant* 254
 8-4b *The Heteroskedasticity Function Must Be Estimated: Feasible GLS* 259

8-4c *What If the Assumed Heteroskedasticity Function Is Wrong?* 262
 8-4d *Prediction and Prediction Intervals with Heteroskedasticity* 264

8-5 The Linear Probability Model Revisited 265

Summary 267
 Key Terms 268
 Problems 268
 Computer Exercises 270

CHAPTER 9 More on Specification and Data Issues 274

9-1 Functional Form Misspecification 275
 9-1a *RESET as a General Test for Functional Form Misspecification* 277
 9-1b *Tests against Nonnested Alternatives* 278
 9-2 Using Proxy Variables for Unobserved Explanatory Variables 279
 9-2a *Using Lagged Dependent Variables as Proxy Variables* 283
 9-2b *A Different Slant on Multiple Regression* 284
 9-3 Models with Random Slopes 285
 9-4 Properties of OLS under Measurement Error 287
 9-4a *Measurement Error in the Dependent Variable* 287
 9-4b *Measurement Error in an Explanatory Variable* 289
 9-5 Missing Data, Nonrandom Samples, and Outlying Observations 293
 9-5a *Missing Data* 293
 9-5b *Nonrandom Samples* 294
 9-5c *Outliers and Influential Observations* 296
 9-6 Least Absolute Deviations Estimation 300
 Summary 302
 Key Terms 303
 Problems 303
 Computer Exercises 307

PART 2

Regression Analysis with Time Series Data 311

CHAPTER 10 Basic Regression Analysis with Time Series Data 312

10-1 The Nature of Time Series Data 312
 10-2 Examples of Time Series Regression Models 313

- 10-2a *Static Models* 314
- 10-2b *Finite Distributed Lag Models* 314
- 10-2c *A Convention about the Time Index* 316

10-3 Finite Sample Properties of OLS under Classical Assumptions 317

- 10-3a *Unbiasedness of OLS* 317
- 10-3b *The Variances of the OLS Estimators and the Gauss-Markov Theorem* 320
- 10-3c *Inference under the Classical Linear Model Assumptions* 322

10-4 Functional Form, Dummy Variables, and Index Numbers 323

10-5 Trends and Seasonality 329

- 10-5a *Characterizing Trending Time Series* 329
- 10-5b *Using Trending Variables in Regression Analysis* 332
- 10-5c *A Detrending Interpretation of Regressions with a Time Trend* 334
- 10-5d *Computing R-Squared When the Dependent Variable Is Trending* 335
- 10-5e *Seasonality* 336

Summary 338

Key Terms 339

Problems 339

Computer Exercises 341

CHAPTER 11 Further Issues in Using OLS with Time Series Data 344

11-1 Stationary and Weakly Dependent Time Series 345

- 11-1a *Stationary and Nonstationary Time Series* 345
- 11-1b *Weakly Dependent Time Series* 346

11-2 Asymptotic Properties of OLS 348

11-3 Using Highly Persistent Time Series in Regression Analysis 354

- 11-3a *Highly Persistent Time Series* 354
- 11-3b *Transformations on Highly Persistent Time Series* 358
- 11-3c *Deciding Whether a Time Series Is $I(1)$* 359

11-4 Dynamically Complete Models and the Absence of Serial Correlation 360

11-5 The Homoskedasticity Assumption for Time Series Models 363

Summary 364

Key Terms 365

Problems 365

Computer Exercises 368

CHAPTER 12 Serial Correlation and Heteroskedasticity in Time Series Regressions 372

12-1 Properties of OLS with Serially Correlated Errors 373

- 12-1a *Unbiasedness and Consistency* 373
- 12-1b *Efficiency and Inference* 373
- 12-1c *Goodness of Fit* 374
- 12-1d *Serial Correlation in the Presence of Lagged Dependent Variables* 374

12-2 Testing for Serial Correlation 376

- 12-2a *A t Test for $AR(1)$ Serial Correlation with Strictly Exogenous Regressors* 376
- 12-2b *The Durbin-Watson Test under Classical Assumptions* 378
- 12-2c *Testing for $AR(1)$ Serial Correlation without Strictly Exogenous Regressors* 379
- 12-2d *Testing for Higher Order Serial Correlation* 380

12-3 Correcting for Serial Correlation with Strictly Exogenous Regressors 381

- 12-3a *Obtaining the Best Linear Unbiased Estimator in the $AR(1)$ Model* 382
- 12-3b *Feasible GLS Estimation with $AR(1)$ Errors* 383
- 12-3c *Comparing OLS and FGLS* 385
- 12-3d *Correcting for Higher Order Serial Correlation* 386

12-4 Differencing and Serial Correlation 387

12-5 Serial Correlation–Robust Inference after OLS 388

12-6 Heteroskedasticity in Time Series Regressions 391

- 12-6a *Heteroskedasticity-Robust Statistics* 392
- 12-6b *Testing for Heteroskedasticity* 392
- 12-6c *Autoregressive Conditional Heteroskedasticity* 393
- 12-6d *Heteroskedasticity and Serial Correlation in Regression Models* 395

Summary 396

Key Terms 396

Problems 396

Computer Exercises 397

PART 3**Advanced Topics 401****CHAPTER 13 Pooling Cross Sections across Time: Simple Panel Data Methods 402**

- 13-1** Pooling Independent Cross Sections across Time 403
 - 13-1a *The Chow Test for Structural Change across Time* 407
- 13-2** Policy Analysis with Pooled Cross Sections 407
- 13-3** Two-Period Panel Data Analysis 412
 - 13-3a *Organizing Panel Data* 417
- 13-4** Policy Analysis with Two-Period Panel Data 417
- 13-5** Differencing with More Than Two Time Periods 420
 - 13-5a *Potential Pitfalls in First Differencing Panel Data* 424
- Summary 424
- Key Terms 425
- Problems 425
- Computer Exercises 426
- Appendix 13A 432

CHAPTER 14 Advanced Panel Data Methods 434

- 14-1** Fixed Effects Estimation 435
 - 14-1a *The Dummy Variable Regression* 438
 - 14-1b *Fixed Effects or First Differencing?* 439
 - 14-1c *Fixed Effects with Unbalanced Panels* 440
- 14-2** Random Effects Models 441
 - 14-2a *Random Effects or Fixed Effects?* 444
- 14-3** The Correlated Random Effects Approach 445
 - 14-3a *Unbalanced Panels* 447
- 14-4** Applying Panel Data Methods to Other Data Structures 448
- Summary 450
- Key Terms 451
- Problems 451
- Computer Exercises 453
- Appendix 14A 457

CHAPTER 15 Instrumental Variables Estimation and Two Stage Least Squares 461

- 15-1** Motivation: Omitted Variables in a Simple Regression Model 462
 - 15-1a *Statistical Inference with the IV Estimator* 466
 - 15-1b *Properties of IV with a Poor Instrumental Variable* 469
 - 15-1c *Computing R-Squared after IV Estimation* 471
- 15-2** IV Estimation of the Multiple Regression Model 471
- 15-3** Two Stage Least Squares 475
 - 15-3a *A Single Endogenous Explanatory Variable* 475
 - 15-3b *Multicollinearity and 2SLS* 477
 - 15-3c *Detecting Weak Instruments* 478
 - 15-3d *Multiple Endogenous Explanatory Variables* 478
 - 15-3e *Testing Multiple Hypotheses after 2SLS Estimation* 479
- 15-4** IV Solutions to Errors-in-Variables Problems 479
- 15-5** Testing for Endogeneity and Testing Overidentifying Restrictions 481
 - 15-5a *Testing for Endogeneity* 481
 - 15-5b *Testing Overidentification Restrictions* 482
- 15-6** 2SLS with Heteroskedasticity 484
- 15-7** Applying 2SLS to Time Series Equations 485
- 15-8** Applying 2SLS to Pooled Cross Sections and Panel Data 487
- Summary 488
- Key Terms 489
- Problems 489
- Computer Exercises 492
- Appendix 15A 496

CHAPTER 16 Simultaneous Equations Models 499

- 16-1** The Nature of Simultaneous Equations Models 500
- 16-2** Simultaneity Bias in OLS 503
- 16-3** Identifying and Estimating a Structural Equation 504
 - 16-3a *Identification in a Two-Equation System* 505
 - 16-3b *Estimation by 2SLS* 508
- 16-4** Systems with More Than Two Equations 510
 - 16-4a *Identification in Systems with Three or More Equations* 510
 - 16-4b *Estimation* 511

- 16-5 Simultaneous Equations Models with Time Series 511
- 16-6 Simultaneous Equations Models with Panel Data 514
 - Summary 516
 - Key Terms 517
 - Problems 517
 - Computer Exercises 519

CHAPTER 17 Limited Dependent Variable Models and Sample Selection Corrections 524

- 17-1 Logit and Probit Models for Binary Response 525
 - 17-1a *Specifying Logit and Probit Models* 525
 - 17-1b *Maximum Likelihood Estimation of Logit and Probit Models* 528
 - 17-1c *Testing Multiple Hypotheses* 529
 - 17-1d *Interpreting the Logit and Probit Estimates* 530
- 17-2 The Tobit Model for Corner Solution Responses 536
 - 17-2a *Interpreting the Tobit Estimates* 537
 - 17-2b *Specification Issues in Tobit Models* 543
- 17-3 The Poisson Regression Model 543
- 17-4 Censored and Truncated Regression Models 547
 - 17-4a *Censored Regression Models* 548
 - 17-4b *Truncated Regression Models* 551
- 17-5 Sample Selection Corrections 553
 - 17-5a *When Is OLS on the Selected Sample Consistent?* 553
 - 17-5b *Incidental Truncation* 554
- Summary 558
- Key Terms 558
- Problems 559
- Computer Exercises 560
- Appendix 17A 565
- Appendix 17B 566

CHAPTER 18 Advanced Time Series Topics 568

- 18-1 Infinite Distributed Lag Models 569
 - 18-1a *The Geometric (or Koyck) Distributed Lag* 571
 - 18-1b *Rational Distributed Lag Models* 572
- 18-2 Testing for Unit Roots 574
- 18-3 Spurious Regression 578
- 18-4 Cointegration and Error Correction Models 580
 - 18-4a *Cointegration* 580
 - 18-4b *Error Correction Models* 584
- 18-5 Forecasting 586

- 18-5a *Types of Regression Models Used for Forecasting* 587
- 18-5b *One-Step-Ahead Forecasting* 588
- 18-5c *Comparing One-Step-Ahead Forecasts* 591
- 18-5d *Multiple-Step-Ahead Forecasts* 592
- 18-5e *Forecasting Trending, Seasonal, and Integrated Processes* 594

- Summary 598
- Key Terms 599
- Problems 600
- Computer Exercises 601

CHAPTER 19 Carrying Out an Empirical Project 605

- 19-1 Posing a Question 605
- 19-2 Literature Review 607
- 19-3 Data Collection 608
 - 19-3a *Deciding on the Appropriate Data Set* 608
 - 19-3b *Entering and Storing Your Data* 609
 - 19-3c *Inspecting, Cleaning, and Summarizing Your Data* 610
- 19-4 Econometric Analysis 611
- 19-5 Writing an Empirical Paper 614
 - 19-5a *Introduction* 614
 - 19-5b *Conceptual (or Theoretical) Framework* 615
 - 19-5c *Econometric Models and Estimation Methods* 615
 - 19-5d *The Data* 617
 - 19-5e *Results* 618
 - 19-5f *Conclusions* 618
 - 19-5g *Style Hints* 619
- Summary 621
- Key Terms 621
- Sample Empirical Projects 621
- List of Journals 626
- Data Sources 627

APPENDIX A Basic Mathematical Tools 628

- A-1 The Summation Operator and Descriptive Statistics 628
- A-2 Properties of Linear Functions 630
- A-3 Proportions and Percentages 633
- A-4 Some Special Functions and their Properties 634
 - A-4a *Quadratic Functions* 634
 - A-4b *The Natural Logarithm* 636
 - A-4c *The Exponential Function* 639

A-5 Differential Calculus 640

- Summary 642
- Key Terms 642
- Problems 643

APPENDIX B Fundamentals of Probability 645

B-1 Random Variables and Their Probability Distributions 645

- B-1a *Discrete Random Variables* 646
- B-1b *Continuous Random Variables* 648

B-2 Joint Distributions, Conditional Distributions, and Independence 649

- B-2a *Joint Distributions and Independence* 649
- B-2b *Conditional Distributions* 651

B-3 Features of Probability Distributions 652

- B-3a *A Measure of Central Tendency: The Expected Value* 652
- B-3b *Properties of Expected Values* 653
- B-3c *Another Measure of Central Tendency: The Median* 655
- B-3d *Measures of Variability: Variance and Standard Deviation* 656
- B-3e *Variance* 656
- B-3f *Standard Deviation* 657
- B-3g *Standardizing a Random Variable* 657
- B-3h *Skewness and Kurtosis* 658

B-4 Features of Joint and Conditional Distributions 658

- B-4a *Measures of Association: Covariance and Correlation* 658
- B-4b *Covariance* 658
- B-4c *Correlation Coefficient* 659
- B-4d *Variance of Sums of Random Variables* 660
- B-4e *Conditional Expectation* 661
- B-4f *Properties of Conditional Expectation* 663
- B-4g *Conditional Variance* 665

B-5 The Normal and Related Distributions 665

- B-5a *The Normal Distribution* 665
- B-5b *The Standard Normal Distribution* 666
- B-5c *Additional Properties of the Normal Distribution* 668
- B-5d *The Chi-Square Distribution* 669
- B-5e *The *t* Distribution* 669
- B-5f *The *F* Distribution* 670

- Summary 672
- Key Terms 672
- Problems 672

APPENDIX C Fundamentals of Mathematical Statistics 674

C-1 Populations, Parameters, and Random Sampling 674

- C-1a *Sampling* 674

C-2 Finite Sample Properties of Estimators 675

- C-2a *Estimators and Estimates* 675
- C-2b *Unbiasedness* 676
- C-2d *The Sampling Variance of Estimators* 678
- C-2e *Efficiency* 679

C-3 Asymptotic or Large Sample Properties of Estimators 681

- C-3a *Consistency* 681
- C-3b *Asymptotic Normality* 683

C-4 General Approaches to Parameter Estimation 684

- C-4a *Method of Moments* 685
- C-4b *Maximum Likelihood* 685
- C-4c *Least Squares* 686

C-5 Interval Estimation and Confidence Intervals 687

- C-5a *The Nature of Interval Estimation* 687
- C-5b *Confidence Intervals for the Mean from a Normally Distributed Population* 689
- C-5c *A Simple Rule of Thumb for a 95% Confidence Interval* 691
- C-5d *Asymptotic Confidence Intervals for Nonnormal Populations* 692

C.6 Hypothesis Testing 693

- C.6a *Fundamentals of Hypothesis Testing* 693
- C.6b *Testing Hypotheses about the Mean in a Normal Population* 695
- C.6c *Asymptotic Tests for Nonnormal Populations* 698
- C.6d *Computing and Using *p*-Values* 698
- C.6e *The Relationship between Confidence Intervals and Hypothesis Testing* 701
- C.6f *Practical versus Statistical Significance* 702

C.7 Remarks on Notation 703

- Summary 703
- Key Terms 704
- Problems 704

APPENDIX D Summary of Matrix Algebra 709

D-1 Basic Definitions 709

D-2 Matrix Operations 710

- D-2a *Matrix Addition* 710

D-2b	<i>Scalar Multiplication</i>	710
D-2c	<i>Matrix Multiplication</i>	711
D-2d	<i>Transpose</i>	712
D-2e	<i>Partitioned Matrix Multiplication</i>	712
D-2f	<i>Trace</i>	713
D-2g	<i>Inverse</i>	713
D-3	Linear Independence and Rank of a Matrix	714
D-4	Quadratic Forms and Positive Definite Matrices	714
D-5	Idempotent Matrices	715
D-6	Differentiation of Linear and Quadratic Forms	715
D-7	Moments and Distributions of Random Vectors	716
D-7a	<i>Expected Value</i>	716
D-7b	<i>Variance-Covariance Matrix</i>	716
D-7c	<i>Multivariate Normal Distribution</i>	716
D-7d	<i>Chi-Square Distribution</i>	717
D-7e	<i>t Distribution</i>	717
D-7f	<i>F Distribution</i>	717
	Summary	717
	Key Terms	717
	Problems	718

APPENDIX E The Linear Regression Model in Matrix Form 720

E-1	The Model and Ordinary Least Squares Estimation	720
E-1a	<i>The Frisch-Waugh Theorem</i>	722
E-2	Finite Sample Properties of OLS	723
E-3	Statistical Inference	726
E-4	Some Asymptotic Analysis	728
E-4a	<i>Wald Statistics for Testing Multiple Hypotheses</i>	730
	Summary	731
	Key Terms	731
	Problems	731

APPENDIX F Answers to Chapter Questions 734

APPENDIX G Statistical Tables 743

References	750
Glossary	756
Index	771

Preface

My motivation for writing the first edition of *Introductory Econometrics: A Modern Approach* was that I saw a fairly wide gap between how econometrics is taught to undergraduates and how empirical researchers think about and apply econometric methods. I became convinced that teaching introductory econometrics from the perspective of professional users of econometrics would actually simplify the presentation, in addition to making the subject much more interesting.

Based on the positive reactions to earlier editions, it appears that my hunch was correct. Many instructors, having a variety of backgrounds and interests and teaching students with different levels of preparation, have embraced the modern approach to econometrics espoused in this text. The emphasis in this edition is still on applying econometrics to real-world problems. Each econometric method is motivated by a particular issue facing researchers analyzing nonexperimental data. The focus in the main text is on understanding and interpreting the assumptions in light of actual empirical applications: the mathematics required is no more than college algebra and basic probability and statistics.

Organized for Today's Econometrics Instructor

The sixth edition preserves the overall organization of the fifth. The most noticeable feature that distinguishes this text from most others is the separation of topics by the kind of data being analyzed. This is a clear departure from the traditional approach, which presents a linear model, lists all assumptions that may be needed at some future point in the analysis, and then proves or asserts results without clearly connecting them to the assumptions. My approach is first to treat, in Part 1, multiple regression analysis with cross-sectional data, under the assumption of random sampling. This setting is natural to students because they are familiar with random sampling from a population in their introductory statistics courses. Importantly, it allows us to distinguish assumptions made about the underlying population regression model—assumptions that can be given economic or behavioral content—from assumptions about how the data were sampled. Discussions about the consequences of nonrandom sampling can be treated in an intuitive fashion after the students have a good grasp of the multiple regression model estimated using random samples.

An important feature of a modern approach is that the explanatory variables—along with the dependent variable—are treated as outcomes of random variables. For the social sciences, allowing random explanatory variables is much more realistic than the traditional assumption of nonrandom explanatory variables. As a nontrivial benefit, the population model/random sampling approach reduces the number of assumptions that students must absorb and understand. Ironically, the classical approach to regression analysis, which treats the explanatory variables as fixed in repeated samples and is still pervasive in introductory texts, literally applies to data collected in an experimental setting. In addition, the contortions required to state and explain assumptions can be confusing to students.

My focus on the population model emphasizes that the fundamental assumptions underlying regression analysis, such as the zero mean assumption on the unobservable error term, are properly

stated conditional on the explanatory variables. This leads to a clear understanding of the kinds of problems, such as heteroskedasticity (nonconstant variance), that can invalidate standard inference procedures. By focusing on the population, I am also able to dispel several misconceptions that arise in econometrics texts at all levels. For example, I explain why the usual R -squared is still valid as a goodness-of-fit measure in the presence of heteroskedasticity (Chapter 8) or serially correlated errors (Chapter 12); I provide a simple demonstration that tests for functional form should not be viewed as general tests of omitted variables (Chapter 9); and I explain why one should always include in a regression model extra control variables that are uncorrelated with the explanatory variable of interest, which is often a key policy variable (Chapter 6).

Because the assumptions for cross-sectional analysis are relatively straightforward yet realistic, students can get involved early with serious cross-sectional applications without having to worry about the thorny issues of trends, seasonality, serial correlation, high persistence, and spurious regression that are ubiquitous in time series regression models. Initially, I figured that my treatment of regression with cross-sectional data followed by regression with time series data would find favor with instructors whose own research interests are in applied microeconomics, and that appears to be the case. It has been gratifying that adopters of the text with an applied time series bent have been equally enthusiastic about the structure of the text. By postponing the econometric analysis of time series data, I am able to put proper focus on the potential pitfalls in analyzing time series data that do not arise with cross-sectional data. In effect, time series econometrics finally gets the serious treatment it deserves in an introductory text.

As in the earlier editions, I have consciously chosen topics that are important for reading journal articles and for conducting basic empirical research. Within each topic, I have deliberately omitted many tests and estimation procedures that, while traditionally included in textbooks, have not withstood the empirical test of time. Likewise, I have emphasized more recent topics that have clearly demonstrated their usefulness, such as obtaining test statistics that are robust to heteroskedasticity (or serial correlation) of unknown form, using multiple years of data for policy analysis, or solving the omitted variable problem by instrumental variables methods. I appear to have made fairly good choices, as I have received only a handful of suggestions for adding or deleting material.

I take a systematic approach throughout the text, by which I mean that each topic is presented by building on the previous material in a logical fashion, and assumptions are introduced only as they are needed to obtain a conclusion. For example, empirical researchers who use econometrics in their research understand that not all of the Gauss-Markov assumptions are needed to show that the ordinary least squares (OLS) estimators are unbiased. Yet the vast majority of econometrics texts introduce a complete set of assumptions (many of which are redundant or in some cases even logically conflicting) before proving the unbiasedness of OLS. Similarly, the normality assumption is often included among the assumptions that are needed for the Gauss-Markov Theorem, even though it is fairly well known that normality plays no role in showing that the OLS estimators are the best linear unbiased estimators.

My systematic approach is illustrated by the order of assumptions that I use for multiple regression in Part 1. This structure results in a natural progression for briefly summarizing the role of each assumption:

MLR.1: Introduce the population model and interpret the population parameters (which we hope to estimate).

MLR.2: Introduce random sampling from the population and describe the data that we use to estimate the population parameters.

MLR.3: Add the assumption on the explanatory variables that allows us to compute the estimates from our sample; this is the so-called no perfect collinearity assumption.

MLR.4: Assume that, in the population, the mean of the unobservable error does not depend on the values of the explanatory variables; this is the “mean independence” assumption combined with a zero population mean for the error, and it is the key assumption that delivers unbiasedness of OLS.

After introducing Assumptions MLR.1 to MLR.3, one can discuss the algebraic properties of ordinary least squares—that is, the properties of OLS for a particular set of data. By adding Assumption MLR.4, we can show that OLS is unbiased (and consistent). Assumption MLR.5 (homoskedasticity) is added for the Gauss-Markov Theorem and for the usual OLS variance formulas to be valid. Assumption MLR.6 (normality), which is not introduced until Chapter 4, is added to round out the classical linear model assumptions. The six assumptions are used to obtain exact statistical inference and to conclude that the OLS estimators have the smallest variances among all unbiased estimators.

I use parallel approaches when I turn to the study of large-sample properties and when I treat regression for time series data in Part 2. The careful presentation and discussion of assumptions makes it relatively easy to transition to Part 3, which covers advanced topics that include using pooled cross-sectional data, exploiting panel data structures, and applying instrumental variables methods. Generally, I have strived to provide a unified view of econometrics, where all estimators and test statistics are obtained using just a few intuitively reasonable principles of estimation and testing (which, of course, also have rigorous justification). For example, regression-based tests for heteroskedasticity and serial correlation are easy for students to grasp because they already have a solid understanding of regression. This is in contrast to treatments that give a set of disjointed recipes for outdated econometric testing procedures.

Throughout the text, I emphasize *ceteris paribus* relationships, which is why, after one chapter on the simple regression model, I move to multiple regression analysis. The multiple regression setting motivates students to think about serious applications early. I also give prominence to policy analysis with all kinds of data structures. Practical topics, such as using proxy variables to obtain *ceteris paribus* effects and interpreting partial effects in models with interaction terms, are covered in a simple fashion.

New to This Edition

I have added new exercises to almost every chapter, including the appendices. Most of the new computer exercises use new data sets, including a data set on student performance and attending a Catholic high school and a time series data set on presidential approval ratings and gasoline prices. I have also added some harder problems that require derivations.

There are several changes to the text worth noting. Chapter 2 contains a more extensive discussion about the relationship between the simple regression coefficient and the correlation coefficient. Chapter 3 clarifies issues with comparing R-squareds from models when data are missing on some variables (thereby reducing sample sizes available for regressions with more explanatory variables).

Chapter 6 introduces the notion of an average partial effect (APE) for models linear in the parameters but including nonlinear functions, primarily quadratics and interaction terms. The notion of an APE, which was implicit in previous editions, has become an important concept in empirical work; understanding how to compute and interpret APEs in the context of OLS is a valuable skill. For more advanced classes, the introduction in Chapter 6 eases the way to the discussion of APEs in the nonlinear models studied in Chapter 17, which also includes an expanded discussion of APEs—including now showing APEs in tables alongside coefficients in logit, probit, and Tobit applications.

In Chapter 8, I refine some of the discussion involving the issue of heteroskedasticity, including an expanded discussion of Chow tests and a more precise description of weighted least squares when the weights must be estimated. Chapter 9, which contains some optional, slightly more advanced topics, defines terms that appear often in the large literature on missing data. A common practice in empirical work is to create indicator variables for missing data, and to include them in a multiple regression analysis. Chapter 9 discusses how this method can be implemented and when it will produce unbiased and consistent estimators.

The treatment of unobserved effects panel data models in chapter 14 has been expanded to include more of a discussion of unbalanced panel data sets, including how the fixed effects, random effects, and correlated random effects approaches still can be applied. Another important addition is a much more detailed discussion on applying fixed effects and random effects methods to cluster samples. I also include discussion of some subtle issues that can arise in using clustered standard errors when the data have been obtained from a random sampling scheme.

Chapter 15 now has a more detailed discussion of the problem of weak instrumental variables so that students can access the basics without having to track down more advanced sources.

Targeted at Undergraduates, Adaptable for Master's Students

The text is designed for undergraduate economics majors who have taken college algebra and one semester of introductory probability and statistics. (Appendices A, B, and C contain the requisite background material.) A one-semester or one-quarter econometrics course would not be expected to cover all, or even any, of the more advanced material in Part 3. A typical introductory course includes Chapters 1 through 8, which cover the basics of simple and multiple regression for cross-sectional data. Provided the emphasis is on intuition and interpreting the empirical examples, the material from the first eight chapters should be accessible to undergraduates in most economics departments. Most instructors will also want to cover at least parts of the chapters on regression analysis with time series data, Chapters 10 and 12, in varying degrees of depth. In the one-semester course that I teach at Michigan State, I cover Chapter 10 fairly carefully, give an overview of the material in Chapter 11, and cover the material on serial correlation in Chapter 12. I find that this basic one-semester course puts students on a solid footing to write empirical papers, such as a term paper, a senior seminar paper, or a senior thesis. Chapter 9 contains more specialized topics that arise in analyzing cross-sectional data, including data problems such as outliers and nonrandom sampling; for a one-semester course, it can be skipped without loss of continuity.

The structure of the text makes it ideal for a course with a cross-sectional or policy analysis focus: the time series chapters can be skipped in lieu of topics from Chapters 9 or 15. Chapter 13 is advanced only in the sense that it treats two new data structures: independently pooled cross sections and two-period panel data analysis. Such data structures are especially useful for policy analysis, and the chapter provides several examples. Students with a good grasp of Chapters 1 through 8 will have little difficulty with Chapter 13. Chapter 14 covers more advanced panel data methods and would probably be covered only in a second course. A good way to end a course on cross-sectional methods is to cover the rudiments of instrumental variables estimation in Chapter 15.

I have used selected material in Part 3, including Chapters 13 and 17, in a senior seminar geared to producing a serious research paper. Along with the basic one-semester course, students who have been exposed to basic panel data analysis, instrumental variables estimation, and limited dependent variable models are in a position to read large segments of the applied social sciences literature. Chapter 17 provides an introduction to the most common limited dependent variable models.

The text is also well suited for an introductory master's level course, where the emphasis is on applications rather than on derivations using matrix algebra. Several instructors have used the text to teach policy analysis at the master's level. For instructors wanting to present the material in matrix form, Appendices D and E are self-contained treatments of the matrix algebra and the multiple regression model in matrix form.

At Michigan State, PhD students in many fields that require data analysis—including accounting, agricultural economics, development economics, economics of education, finance, international economics, labor economics, macroeconomics, political science, and public finance—have found the text

to be a useful bridge between the empirical work that they read and the more theoretical econometrics they learn at the PhD level.

Design Features

Numerous in-text questions are scattered throughout, with answers supplied in Appendix F. These questions are intended to provide students with immediate feedback. Each chapter contains many numbered examples. Several of these are case studies drawn from recently published papers, but where I have used my judgment to simplify the analysis, hopefully without sacrificing the main point. The end-of-chapter problems and computer exercises are heavily oriented toward empirical work, rather than complicated derivations. The students are asked to reason carefully based on what they have learned. The computer exercises often expand on the in-text examples. Several exercises use data sets from published works or similar data sets that are motivated by published research in economics and other fields.

A pioneering feature of this introductory econometrics text is the extensive glossary. The short definitions and descriptions are a helpful refresher for students studying for exams or reading empirical research that uses econometric methods. I have added and updated several entries for the fifth edition.

Data Sets—Available in Six Formats

This edition adds R data set as an additional format for viewing and analyzing data. In response to popular demand, this edition also provides the Minitab[®] format. With more than 100 data sets in six different formats, including Stata[®], EViews[®], Minitab[®], Microsoft[®] Excel, and R, the instructor has many options for problem sets, examples, and term projects. Because most of the data sets come from actual research, some are very large. Except for partial lists of data sets to illustrate the various data structures, the data sets are not reported in the text. This book is geared to a course where computer work plays an integral role.

Updated Data Sets Handbook

An extensive data description manual is also available online. This manual contains a list of data sources along with suggestions for ways to use the data sets that are not described in the text. This unique handbook, created by author Jeffrey M. Wooldridge, lists the source of all data sets for quick reference and how each might be used. Because the data book contains page numbers, it is easy to see how the author used the data in the text. Students may want to view the descriptions of each data set and it can help guide instructors in generating new homework exercises, exam problems, or term projects. The author also provides suggestions on improving the data sets in this detailed resource that is available on the book's companion website at <http://login.cengage.com> and students can access it free at www.cengagebrain.com.

Instructor Supplements

Instructor's Manual with Solutions

The *Instructor's Manual with Solutions* contains answers to all problems and exercises, as well as teaching tips on how to present the material in each chapter. The instructor's manual also contains

sources for each of the data files, with many suggestions for how to use them on problem sets, exams, and term papers. This supplement is available online only to instructors at <http://login.cengage.com>.

PowerPoint Slides

Exceptional PowerPoint® presentation slides help you create engaging, memorable lectures. You will find teaching slides for each chapter in this edition, including the advanced chapters in Part 3. You can modify or customize the slides for your specific course. PowerPoint® slides are available for convenient download on the instructor-only, password-protected portion of the book's companion website at <http://login.cengage.com>.

Scientific Word Slides

Developed by the author, Scientific Word® slides offer an alternative format for instructors who prefer the Scientific Word® platform, the word processor created by MacKichan Software, Inc. for composing mathematical and technical documents using LaTeX typesetting. These slides are based on the author's actual lectures and are available in PDF and TeX formats for convenient download on the instructor-only, password-protected section of the book's companion website at <http://login.cengage.com>.

Test Bank

Cengage Learning Testing, powered by Cognero® is a flexible, online system that allows you to import, edit, and manipulate content from the text's test bank or elsewhere. You have the flexibility to include your own favorite test questions, create multiple test versions in an instant, and deliver tests from your LMS, your classroom, or wherever you want. In the test bank for INTRODUCTORY ECONOMETRICS, 6E you will find a wealth and variety of problems, ranging from multiple-choice to questions that require simple statistical derivations to questions that require interpreting computer output.

Student Supplements

MindTap

MindTap® for INTRODUCTORY ECONOMETRICS, 6E provides you with the tools you need to better manage your limited time—you can complete assignments whenever and wherever you are ready to learn with course material specially customized by your instructor and streamlined in one proven, easy-to-use interface. With an array of tools and apps—from note taking to flashcards—you will get a true understanding of course concepts, helping you to achieve better grades and setting the groundwork for your future courses.

Aplia

Millions of students use Aplia™ to better prepare for class and for their exams. Aplia assignments mean “no surprises”—with an at-a-glance view of current assignments organized by due date. You always know what's due, and when. Aplia ties your lessons into real-world applications so you get a bigger, better picture of how you'll use your education in your future workplace. Automatic grading and immediate feedback helps you master content the right way the first time.

Student Solutions Manual

Now you can maximize your study time and further your course success with this dynamic online resource. This helpful Solutions Manual includes detailed steps and solutions to odd-numbered problems as well as computer exercises in the text. This supplement is available as a free resource at www.cengagebrain.com.

Suggestions for Designing Your Course

I have already commented on the contents of most of the chapters as well as possible outlines for courses. Here I provide more specific comments about material in chapters that might be covered or skipped:

Chapter 9 has some interesting examples (such as a wage regression that includes IQ score as an explanatory variable). The rubric of proxy variables does not have to be formally introduced to present these kinds of examples, and I typically do so when finishing up cross-sectional analysis. In Chapter 12, for a one-semester course, I skip the material on serial correlation robust inference for ordinary least squares as well as dynamic models of heteroskedasticity.

Even in a second course I tend to spend only a little time on Chapter 16, which covers simultaneous equations analysis. I have found that instructors differ widely in their opinions on the importance of teaching simultaneous equations models to undergraduates. Some think this material is fundamental; others think it is rarely applicable. My own view is that simultaneous equations models are overused (see Chapter 16 for a discussion). If one reads applications carefully, omitted variables and measurement error are much more likely to be the reason one adopts instrumental variables estimation, and this is why I use omitted variables to motivate instrumental variables estimation in Chapter 15. Still, simultaneous equations models are indispensable for estimating demand and supply functions, and they apply in some other important cases as well.

Chapter 17 is the only chapter that considers models inherently nonlinear in their parameters, and this puts an extra burden on the student. The first material one should cover in this chapter is on probit and logit models for binary response. My presentation of Tobit models and censored regression still appears to be novel in introductory texts. I explicitly recognize that the Tobit model is applied to corner solution outcomes on random samples, while censored regression is applied when the data collection process censors the dependent variable at essentially arbitrary thresholds.

Chapter 18 covers some recent important topics from time series econometrics, including testing for unit roots and cointegration. I cover this material only in a second-semester course at either the undergraduate or master's level. A fairly detailed introduction to forecasting is also included in Chapter 18.

Chapter 19, which would be added to the syllabus for a course that requires a term paper, is much more extensive than similar chapters in other texts. It summarizes some of the methods appropriate for various kinds of problems and data structures, points out potential pitfalls, explains in some detail how to write a term paper in empirical economics, and includes suggestions for possible projects.

Acknowledgments

I would like to thank those who reviewed and provided helpful comments for this and previous editions of the text:

Erica Johnson, *Gonzaga University*

Mary Ellen Benedict, *Bowling Green State University*

Yan Li, *Temple University*

Melissa Tartari,
Yale University

- Michael Allgrunn, *University of South Dakota*
- Gregory Colman, *Pace University*
- Yoo-Mi Chin, *Missouri University of Science and Technology*
- Arsen Melkumian, *Western Illinois University*
- Kevin J. Murphy, *Oakland University*
- Kristine Grimsrud, *University of New Mexico*
- Will Melick, *Kenyon College*
- Philip H. Brown, *Colby College*
- Argun Saatcioglu, *University of Kansas*
- Ken Brown, *University of Northern Iowa*
- Michael R. Jonas, *University of San Francisco*
- Melissa Yeoh, *Berry College*
- Nikolaos Papanikolaou, *SUNY at New Paltz*
- Konstantin Golyaev, *University of Minnesota*
- Soren Hauge, *Ripon College*
- Kevin Williams, *University of Minnesota*
- Hailong Qian, *Saint Louis University*
- Rod Hissong, *University of Texas at Arlington*
- Steven Cuellar, *Sonoma State University*
- Yanan Di, *Wagner College*
- John Fitzgerald, *Bowdoin College*
- Philip N. Jefferson, *Swarthmore College*
- Yongsheng Wang, *Washington and Jefferson College*
- Sheng-Kai Chang, *National Taiwan University*
- Damayanti Ghosh, *Binghamton University*
- Susan Averett, *Lafayette College*
- Kevin J. Mumford, *Purdue University*
- Nicolai V. Kuminoff, *Arizona State University*
- Subarna K. Samanta, *The College of New Jersey*
- Jing Li, *South Dakota State University*
- Gary Wagner, *University of Arkansas–Little Rock*
- Kelly Cobourn, *Boise State University*
- Timothy Dittmer, *Central Washington University*
- Daniel Fischmar, *Westminster College*
- Subha Mani, *Fordham University*
- John Maluccio, *Middlebury College*
- James Warner, *College of Wooster*
- Christopher Magee, *Bucknell University*
- Andrew Ewing, *Eckerd College*
- Debra Israel, *Indiana State University*
- Jay Goodliffe, *Brigham Young University*
- Stanley R. Thompson, *The Ohio State University*
- Michael Robinson, *Mount Holyoke College*
- Ivan Jeliakov, *University of California, Irvine*
- Heather O’Neill, *Ursinus College*
- Leslie Papke, *Michigan State University*
- Timothy Vogelsang, *Michigan State University*
- Stephen Woodbury, *Michigan State University*

Some of the changes I discussed earlier were driven by comments I received from people on this list, and I continue to mull over other specific suggestions made by one or more reviewers.

Many students and teaching assistants, too numerous to list, have caught mistakes in earlier editions or have suggested rewording some paragraphs. I am grateful to them.

As always, it was a pleasure working with the team at Cengage Learning. Mike Worls, my long-time Product Director, has learned very well how to guide me with a firm yet gentle hand. Chris Rader has quickly mastered the difficult challenges of being the developmental editor of a dense, technical textbook. His careful reading of the manuscript and fine eye for detail have improved this sixth edition considerably.

This book is dedicated to my wife, Leslie Papke, who contributed materially to this edition by writing the initial versions of the Scientific Word slides for the chapters in Part 3; she then used the slides in her public policy course. Our children have contributed, too: Edmund has helped me keep the data handbook current, and Gwenth keeps us entertained with her artistic talents.

Jeffrey M. Wooldridge

About the Author

Jeffrey M. Wooldridge is University Distinguished Professor of Economics at Michigan State University, where he has taught since 1991. From 1986 to 1991, he was an assistant professor of economics at the Massachusetts Institute of Technology. He received his bachelor of arts, with majors in computer science and economics, from the University of California, Berkeley, in 1982, and received his doctorate in economics in 1986 from the University of California, San Diego. He has published more than 60 articles in internationally recognized journals, as well as several book chapters. He is also the author of *Econometric Analysis of Cross Section and Panel Data*, second edition. His awards include an Alfred P. Sloan Research Fellowship, the Plura Scripsit award from *Econometric Theory*, the Sir Richard Stone prize from the *Journal of Applied Econometrics*, and three graduate teacher-of-the-year awards from MIT. He is a fellow of the Econometric Society and of the *Journal of Econometrics*. He is past editor of the *Journal of Business and Economic Statistics*, and past econometrics coeditor of *Economics Letters*. He has served on the editorial boards of *Econometric Theory*, the *Journal of Economic Literature*, the *Journal of Econometrics*, the *Review of Economics and Statistics*, and the *Stata Journal*. He has also acted as an occasional econometrics consultant for Arthur Andersen, Charles River Associates, the Washington State Institute for Public Policy, Stratus Consulting, and Industrial Economics, Incorporated.

The Nature of Econometrics and Economic Data

Chapter 1 discusses the scope of econometrics and raises general issues that arise in the application of econometric methods. Section 1-1 provides a brief discussion about the purpose and scope of econometrics and how it fits into economic analysis. Section 1-2 provides examples of how one can start with an economic theory and build a model that can be estimated using data. Section 1-3 examines the kinds of data sets that are used in business, economics, and other social sciences. Section 1-4 provides an intuitive discussion of the difficulties associated with the inference of causality in the social sciences.

1-1 What Is Econometrics?

Imagine that you are hired by your state government to evaluate the effectiveness of a publicly funded job training program. Suppose this program teaches workers various ways to use computers in the manufacturing process. The 20-week program offers courses during nonworking hours. Any hourly manufacturing worker may participate, and enrollment in all or part of the program is voluntary. You are to determine what, if any, effect the training program has on each worker's subsequent hourly wage.

Now, suppose you work for an investment bank. You are to study the returns on different investment strategies involving short-term U.S. treasury bills to decide whether they comply with implied economic theories.

The task of answering such questions may seem daunting at first. At this point, you may only have a vague idea of the kind of data you would need to collect. By the end of this introductory econometrics course, you should know how to use econometric methods to formally evaluate a job training program or to test a simple economic theory.

Econometrics is based upon the development of statistical methods for estimating economic relationships, testing economic theories, and evaluating and implementing government and business policy. The most common application of econometrics is the forecasting of such important macroeconomic variables as interest rates, inflation rates, and gross domestic product (GDP). Whereas forecasts of economic indicators are highly visible and often widely published, econometric methods can be used in economic areas that have nothing to do with macroeconomic forecasting. For example, we will study the effects of political campaign expenditures on voting outcomes. We will consider the effect of school spending on student performance in the field of education. In addition, we will learn how to use econometric methods for forecasting economic time series.

Econometrics has evolved as a separate discipline from mathematical statistics because the former focuses on the problems inherent in collecting and analyzing nonexperimental economic data. **Nonexperimental data** are not accumulated through controlled experiments on individuals, firms, or segments of the economy. (Nonexperimental data are sometimes called **observational data**, or **retrospective data**, to emphasize the fact that the researcher is a passive collector of the data.) **Experimental data** are often collected in laboratory environments in the natural sciences, but they are much more difficult to obtain in the social sciences. Although some social experiments can be devised, it is often impossible, prohibitively expensive, or morally repugnant to conduct the kinds of controlled experiments that would be needed to address economic issues. We give some specific examples of the differences between experimental and nonexperimental data in Section 1-4.

Naturally, econometricians have borrowed from mathematical statisticians whenever possible. The method of multiple regression analysis is the mainstay in both fields, but its focus and interpretation can differ markedly. In addition, economists have devised new techniques to deal with the complexities of economic data and to test the predictions of economic theories.

1-2 Steps in Empirical Economic Analysis

Econometric methods are relevant in virtually every branch of applied economics. They come into play either when we have an economic theory to test or when we have a relationship in mind that has some importance for business decisions or policy analysis. An **empirical analysis** uses data to test a theory or to estimate a relationship.

How does one go about structuring an empirical economic analysis? It may seem obvious, but it is worth emphasizing that the first step in any empirical analysis is the careful formulation of the question of interest. The question might deal with testing a certain aspect of an economic theory, or it might pertain to testing the effects of a government policy. In principle, econometric methods can be used to answer a wide range of questions.

In some cases, especially those that involve the testing of economic theories, a formal **economic model** is constructed. An economic model consists of mathematical equations that describe various relationships. Economists are well known for their building of models to describe a vast array of behaviors. For example, in intermediate microeconomics, individual consumption decisions, subject to a budget constraint, are described by mathematical models. The basic premise underlying these models is *utility maximization*. The assumption that individuals make choices to maximize their well-being, subject to resource constraints, gives us a very powerful framework for creating tractable economic models and making clear predictions. In the context of consumption decisions, utility maximization leads to a set of *demand equations*. In a demand equation, the quantity demanded of each commodity depends on the price of the goods, the price of substitute and complementary goods, the consumer's income, and the individual's characteristics that affect taste. These equations can form the basis of an econometric analysis of consumer demand.

Economists have used basic economic tools, such as the utility maximization framework, to explain behaviors that at first glance may appear to be noneconomic in nature. A classic example is Becker's (1968) economic model of criminal behavior.

EXAMPLE 1.1 Economic Model of Crime

In a seminal article, Nobel Prize winner Gary Becker postulated a utility maximization framework to describe an individual's participation in crime. Certain crimes have clear economic rewards, but most criminal behaviors have costs. The opportunity costs of crime prevent the criminal from participating in other activities such as legal employment. In addition, there are costs associated with the possibility of being caught and then, if convicted, the costs associated with incarceration. From Becker's perspective, the decision to undertake illegal activity is one of resource allocation, with the benefits and costs of competing activities taken into account.

Under general assumptions, we can derive an equation describing the amount of time spent in criminal activity as a function of various factors. We might represent such a function as

$$y = f(x_1, x_2, x_3, x_4, x_5, x_6, x_7), \quad [1.1]$$

where

- y = hours spent in criminal activities,
- x_1 = "wage" for an hour spent in criminal activity,
- x_2 = hourly wage in legal employment,
- x_3 = income other than from crime or employment,
- x_4 = probability of getting caught,
- x_5 = probability of being convicted if caught,
- x_6 = expected sentence if convicted, and
- x_7 = age.

Other factors generally affect a person's decision to participate in crime, but the list above is representative of what might result from a formal economic analysis. As is common in economic theory, we have not been specific about the function $f(\cdot)$ in (1.1). This function depends on an underlying utility function, which is rarely known. Nevertheless, we can use economic theory—or introspection—to predict the effect that each variable would have on criminal activity. This is the basis for an econometric analysis of individual criminal activity.

Formal economic modeling is sometimes the starting point for empirical analysis, but it is more common to use economic theory less formally, or even to rely entirely on intuition. You may agree that the determinants of criminal behavior appearing in equation (1.1) are reasonable based on common sense; we might arrive at such an equation directly, without starting from utility maximization. This view has some merit, although there are cases in which formal derivations provide insights that intuition can overlook.

Next is an example of an equation that we can derive through somewhat informal reasoning.

EXAMPLE 1.2 Job Training and Worker Productivity

Consider the problem posed at the beginning of Section 1-1. A labor economist would like to examine the effects of job training on worker productivity. In this case, there is little need for formal economic theory. Basic economic understanding is sufficient for realizing that factors such as education, experience, and training affect worker productivity. Also, economists are well aware that workers are paid commensurate with their productivity. This simple reasoning leads to a model such as

$$wage = f(educ, exper, training), \quad [1.2]$$

where

- $wage$ = hourly wage,
- $educ$ = years of formal education,
- $exper$ = years of workforce experience, and
- $training$ = weeks spent in job training.

Again, other factors generally affect the wage rate, but equation (1.2) captures the essence of the problem.

After we specify an economic model, we need to turn it into what we call an **econometric model**. Because we will deal with econometric models throughout this text, it is important to know how an econometric model relates to an economic model. Take equation (1.1) as an example. The form of the function $f(\cdot)$ must be specified before we can undertake an econometric analysis. A second issue concerning (1.1) is how to deal with variables that cannot reasonably be observed. For example, consider the wage that a person can earn in criminal activity. In principle, such a quantity is well defined, but it would be difficult if not impossible to observe this wage for a given individual. Even variables such as the probability of being arrested cannot realistically be obtained for a given individual, but at least we can observe relevant arrest statistics and derive a variable that approximates the probability of arrest. Many other factors affect criminal behavior that we cannot even list, let alone observe, but we must somehow account for them.

The ambiguities inherent in the economic model of crime are resolved by specifying a particular econometric model:

$$\begin{aligned} crime = & \beta_0 + \beta_1 wage_m + \beta_2 othinc + \beta_3 freqarr + \beta_4 freqconv \\ & + \beta_5 avgsen + \beta_6 age + u, \end{aligned} \quad [1.3]$$

where

- $crime$ = some measure of the frequency of criminal activity,
- $wage_m$ = the wage that can be earned in legal employment,
- $othinc$ = the income from other sources (assets, inheritance, and so on),
- $freqarr$ = the frequency of arrests for prior infractions (to approximate the probability of arrest),
- $freqconv$ = the frequency of conviction, and
- $avgsen$ = the average sentence length after conviction.

The choice of these variables is determined by the economic theory as well as data considerations. The term u contains unobserved factors, such as the wage for criminal activity, moral character, family background, and errors in measuring things like criminal activity and the probability of arrest. We could add family background variables to the model, such as number of siblings, parents' education, and so on, but we can never eliminate u entirely. In fact, dealing with this *error term* or *disturbance term* is perhaps the most important component of any econometric analysis.

The constants $\beta_0, \beta_1, \dots, \beta_6$ are the *parameters* of the econometric model, and they describe the directions and strengths of the relationship between *crime* and the factors used to determine *crime* in the model.

A complete econometric model for Example 1.2 might be

$$wage = \beta_0 + \beta_1 educ + \beta_2 exper + \beta_3 training + u, \quad [1.4]$$

where the term u contains factors such as “innate ability,” quality of education, family background, and the myriad other factors that can influence a person's wage. If we are specifically concerned about the effects of job training, then β_3 is the parameter of interest.

For the most part, econometric analysis begins by specifying an econometric model, without consideration of the details of the model's creation. We generally follow this approach, largely because careful derivation of something like the economic model of crime is time consuming and can take us into some specialized and often difficult areas of economic theory. Economic reasoning will play a role in our examples, and we will merge any underlying economic theory into the econometric model specification. In the economic model of crime example, we would start with an econometric model such as (1.3) and use economic reasoning and common sense as guides for choosing the variables. Although this approach loses some of the richness of economic analysis, it is commonly and effectively applied by careful researchers.

Once an econometric model such as (1.3) or (1.4) has been specified, various *hypotheses* of interest can be stated in terms of the unknown parameters. For example, in equation (1.3), we might hypothesize that $wage_m$, the wage that can be earned in legal employment, has no effect on criminal behavior. In the context of this particular econometric model, the hypothesis is equivalent to $\beta_1 = 0$.

An empirical analysis, by definition, requires data. After data on the relevant variables have been collected, econometric methods are used to estimate the parameters in the econometric model and to formally test hypotheses of interest. In some cases, the econometric model is used to make predictions in either the testing of a theory or the study of a policy's impact.

Because data collection is so important in empirical work, Section 1-3 will describe the kinds of data that we are likely to encounter.

1-3 The Structure of Economic Data

Economic data sets come in a variety of types. Whereas some econometric methods can be applied with little or no modification to many different kinds of data sets, the special features of some data sets must be accounted for or should be exploited. We next describe the most important data structures encountered in applied work.

1-3a Cross-Sectional Data

A **cross-sectional data set** consists of a sample of individuals, households, firms, cities, states, countries, or a variety of other units, taken at a given point in time. Sometimes, the data on all units do not correspond to precisely the same time period. For example, several families may be surveyed during different weeks within a year. In a pure cross-sectional analysis, we would ignore any minor timing differences in collecting the data. If a set of families was surveyed during different weeks of the same year, we would still view this as a cross-sectional data set.

An important feature of cross-sectional data is that we can often assume that they have been obtained by **random sampling** from the underlying population. For example, if we obtain information on wages, education, experience, and other characteristics by randomly drawing 500 people from the working population, then we have a random sample from the population of all working people. Random sampling is the sampling scheme covered in introductory statistics courses, and it simplifies the analysis of cross-sectional data. A review of random sampling is contained in Appendix C.

Sometimes, random sampling is not appropriate as an assumption for analyzing cross-sectional data. For example, suppose we are interested in studying factors that influence the accumulation of family wealth. We could survey a random sample of families, but some families might refuse to report their wealth. If, for example, wealthier families are less likely to disclose their wealth, then the resulting sample on wealth is not a random sample from the population of all families. This is an illustration of a sample selection problem, an advanced topic that we will discuss in Chapter 17.

Another violation of random sampling occurs when we sample from units that are large relative to the population, particularly geographical units. The potential problem in such cases is that the population is not large enough to reasonably assume the observations are independent draws. For example, if we want to explain new business activity across states as a function of wage rates, energy prices, corporate and property tax rates, services provided, quality of the workforce, and other state characteristics, it is unlikely that business activities in states near one another are independent. It turns out that the econometric methods that we discuss do work in such situations, but they sometimes need to be refined. For the most part, we will ignore the intricacies that arise in analyzing such situations and treat these problems in a random sampling framework, even when it is not technically correct to do so.

Cross-sectional data are widely used in economics and other social sciences. In economics, the analysis of cross-sectional data is closely aligned with the applied microeconomics fields, such as labor economics, state and local public finance, industrial organization, urban economics, demography, and health economics. Data on individuals, households, firms, and cities at a given point in time are important for testing microeconomic hypotheses and evaluating economic policies.

The cross-sectional data used for econometric analysis can be represented and stored in computers. Table 1.1 contains, in abbreviated form, a cross-sectional data set on 526 working individuals

obsno	wage	educ	exper	female	married
1	3.10	11	2	1	0
2	3.24	12	22	1	1
3	3.00	11	2	0	0
4	6.00	8	44	0	1
5	5.30	12	7	0	1
.
.
.
525	11.56	16	5	0	1
526	3.50	14	5	1	0

for the year 1976. (This is a subset of the data in the file WAGE1.) The variables include *wage* (in dollars per hour), *educ* (years of education), *exper* (years of potential labor force experience), *female* (an indicator for gender), and *married* (marital status). These last two variables are binary (zero-one) in nature and serve to indicate qualitative features of the individual (the person is female or not; the person is married or not). We will have much to say about binary variables in Chapter 7 and beyond.

The variable *obsno* in Table 1.1 is the observation number assigned to each person in the sample. Unlike the other variables, it is not a characteristic of the individual. All econometrics and statistics software packages assign an observation number to each data unit. Intuition should tell you that, for data such as that in Table 1.1, it does not matter which person is labeled as observation 1, which person is called observation 2, and so on. The fact that the ordering of the data does not matter for econometric analysis is a key feature of cross-sectional data sets obtained from random sampling.

Different variables sometimes correspond to different time periods in cross-sectional data sets. For example, to determine the effects of government policies on long-term economic growth, economists have studied the relationship between growth in real per capita GDP over a certain period (say, 1960 to 1985) and variables determined in part by government policy in 1960 (government consumption as a percentage of GDP and adult secondary education rates). Such a data set might be represented as in Table 1.2, which constitutes part of the data set used in the study of cross-country growth rates by De Long and Summers (1991).

The variable *gpcrgdp* represents average growth in real per capita GDP over the period 1960 to 1985. The fact that *govcons60* (government consumption as a percentage of GDP) and *second60*

obsno	country	gpcrgdp	govcons60	second60
1	Argentina	0.89	9	32
2	Austria	3.32	16	50
3	Belgium	2.56	13	69
4	Bolivia	1.24	18	12
.
.
.
61	Zimbabwe	2.30	17	6